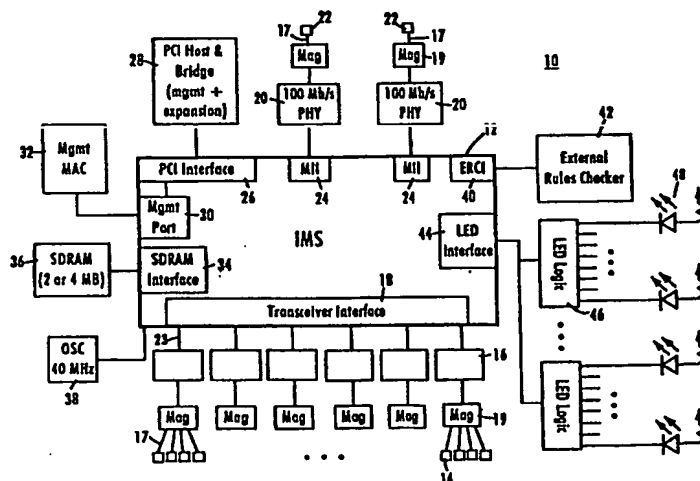




INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification ⁶ : H04Q 3/64	A2	(11) International Publication Number: WO 98/36587 (43) International Publication Date: 20 August 1998 (20.08.98)
(21) International Application Number: PCT/US98/02357 (22) International Filing Date: 6 February 1998 (06.02.98) (30) Priority Data: 60/038,025 14 February 1997 (14.02.97) US 08/993,147 18 December 1997 (18.12.97) US (71) Applicant: ADVANCED MICRO DEVICES, INC. [US/US]; One AMD Place, Mail Stop 68, Sunnyvale, CA 94088-3453 (US). (72) Inventors: ERIMLI, Bahadir; 2101 California Street #109, Mountain View, CA 94040 (US). CRAYFORD, Ian; 5380 Eileen Drive, San Jose, CA 95129 (US). EGBERT, Chandan; 3632 Bloomsbury Way, San Jose, CA 95132 (US). (74) Agent: ZAHRT, William, D., II.; Advanced Micro Devices, Inc., One AMD Place, Mail Stop 68, Sunnyvale, CA 94088-3453 (US).		(81) Designated States: JP, KR, European patent (AT, BE, CH, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE). Published <i>Without international search report and to be republished upon receipt of that report.</i>

(54) Title: QUEUING STRUCTURE AND METHOD FOR PRIORITIZATION OF FRAMES IN A NETWORK SWITCH



(57) Abstract

A method and arrangement for queuing data in a prioritized manner have a queue with a single queue write side in which data entries are input to the queue. The queue also has a plurality of queue read sides, with each read side having a different priority level. An entry exiting from the queue write side is examined to determine the priority level of the entry, and then placed into the queue read side with the matching priority level as the entry. The queue read sides form the output of the queue, and are polled and emptied so that the higher priority queues are emptied completely before emptying lower priority queues.

FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece	ML	Mali	TR	Turkey
BG	Bulgaria	HU	Hungary	MN	Mongolia	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MR	Mauritania	UA	Ukraine
BR	Brazil	IL	Israel	MW	Malawi	UG	Uganda
BY	Belarus	IS	Iceland	MX	Mexico	US	United States of America
CA	Canada	IT	Italy	NE	Niger	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NL	Netherlands	VN	Viet Nam
CG	Congo	KE	Kenya	NO	Norway	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NZ	New Zealand	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	PL	Poland		
CM	Cameroon	KR	Republic of Korea	PT	Portugal		
CN	China	KZ	Kazakstan	RO	Romania		
CU	Cuba	LC	Saint Lucia	RU	Russian Federation		
CZ	Czech Republic	LI	Liechtenstein	SD	Sudan		
DE	Germany	LK	Sri Lanka	SE	Sweden		
DK	Denmark	LR	Liberia	SG	Singapore		
EE	Estonia						

QUEUING STRUCTURE AND METHOD FOR PRIORITIZATION OF FRAMES IN A NETWORK SWITCH

FIELD OF THE INVENTION

The present invention relates to the field of temporary data storage, and more particularly, to queues that receive, temporarily store, and release entries in a prioritized manner.

BACKGROUND ART

In many types of systems, from simple to sophisticated, entries to the system often cannot be processed through the system immediately. Rather than discard the entries, "queues" are employed to temporarily hold the entries and release an entry for processing by the system when the system is able to process the next entry. An example of a queue from everyday life is a customer line, at a bank or an airport ticket counter, for example. Typical queues maintain a first-in, first-out (FIFO) ordering of the entries to the queue.

There are many types of electronic data systems in which queues are used. These include microprocessors, memory transfer systems, airline telephone reservation systems, and packet switched networks, for example. In most systems, it is desirable that the queues have low latencies so that processing of an entry is not delayed very long due to delays caused by the queues themselves. A low queue latency means that an entry will flow from the entrance to the queue to the exit of the queue quickly, in comparison to queues with higher latencies. One factor that has a significant impact on the latency of a queue is the length, or capacity, of the queue. The greater the capacity of the queue to store entries, the higher the latency of the queue.

Many systems carry different types of traffic. Some types of traffic may be considered higher in priority than other types of traffic. For example, video data is typically considered to be a high priority traffic type, since this data must normally be delivered on a high-speed basis to provide an acceptable video signal at the output device. Other types of data may be considered high priority due to their relative importance, so that discarding of such data due to congestion is to be avoided if at all possible. By contrast, low priority frames do not have as rigid time requirements for arrival at their destinations.

A problem that occasionally occurs in network switches and other devices that receive and forward data is internal congestion, which occurs when the device receives more data than it can transmit. For example, a typical network switch includes at least one output queue that queues frames of data received at the switch for transmission from a switch port. If the output queue becomes full, as may happen when data arrives at the switch faster than it can be forwarded, later arriving frames will be discarded, and the data lost. Systems which have an established priority scheme for data attempt to avoid the discarding of the higher priority data. A common arrangement employs multiple queues, with one queue dedicated for each of the different priority levels. Hence, a high priority frame will be loaded into the high priority queue, while a lower priority frame will be loaded into a lower priority queue. In such arrangements, the high priority queue is serviced and emptied before the lower priority queue is serviced.

The problem with such systems is the large amount of space that is devoted to the queuing structure to provide the multiple priorities. For example, at a single port, where there are three different levels of priority allowed per entry to be queued, three separate queues are employed. This requires dedication of usually precious on-chip space to the multiple queues. The problem is especially felt when the switch or device has multiple ports, for example 28 ports. If three separate queues are provided for each of the ports to support three different levels of priority, 84 queues would be located on the chip, increasing the chip size and cost.

SUMMARY OF THE INVENTION

There is a need to provide a queuing structure and a method of queuing that will prioritize traffic in a device in a space and cost-effective manner.

This and other needs are met by embodiments of the present invention which provide a queue that queues and prioritizes entries and has a write side and a plurality of read sides. The queue side forms a queue input at which entries to the queue are received, each entry having a priority level. The plurality of queue read sides forms a queue output at which entries from the queue write side are received, queued and output from the queue. The queue read sides have different levels of queuing priority from a highest priority queue read side to a lowest priority queue read side. Control logic examines entries to the queue write side to determine the priority level of the entries and directs the entries from the queue write side to one of the queue read sides based on the priority levels of the entries.

An advantage of the queue of the present invention is that the entire queue structure is not replicated for each of the different priority levels, and for each of the different ports. Instead, the queue has two sections, a write side and a plurality of read sides. Each of the read sides of the queue can be kept relatively small to maintain the overall capacity of the queue to a desired amount, yet provide multiple levels of priority for queuing different priority data. The single write side, which is relatively much larger than any single one of the read sides, is not duplicated and is used in the queuing of all of the different priority levels of data. Only the relatively small read sides are duplicated to provide the different priority queuing. Hence, the overall size of the queue arrangement for providing different queuing priorities is much less than prior arrangements in which a plurality of entire queues are dedicated to the different priority levels.

The earlier stated needs are also met by another embodiment of the present invention which provides a multiport network switch for a packet switched network, the switch having a plurality of ports through which packets from the network are received and from which packets to the network are transmitted. The switch has a plurality of queues corresponding to the plurality of ports. At least one of those queues has a queue write side forming a queue input at which entries to the queue are received, each entry having a priority level, and a plurality of queue read sides forming a queue output at which entries from the queue write side are received, queued and output from the queue. The plurality of queue read sides have different levels of queuing priority from a highest priority queue read side to a lowest priority queue read side. Control logic is provided that is configured to examine entries to the queue write side to determine the priority level of the entries and to direct the entries from the queue write side to one of the queue read sides based on the priority levels of the entries.

The network switch according to the embodiments of the present invention has the advantage, among others, of providing multiple ports with output queues that will queue data in a prioritized manner. It does so in a

cost-effective and space saving manner due to the duplication of only the read sides of an output queue for the different priority levels, rather than duplication of an entire queue. This split queue structure with multiple read sides is especially advantageous when the switch has multiple ports, since the space savings may be realized at each of the ports.

Another embodiment of the present invention also satisfies the earlier stated needs by providing a method of queuing and prioritizing entries including the steps of receiving entries to the queue in a queue write side that forms an input of the queue and determining a level of priority of the entries to the queue write side. Entries exiting the queue write side are directed to one of a plurality of queue read sides of the queue, based on the determined level of priority of the entries. The queue read sides have different levels of queuing priority, from a highest priority queue read side to a lowest priority queue read side.

The foregoing and other features, aspects and advantages of the present invention will become more apparent from the following detailed description of the present invention when taken in conjunction with the accompanying drawings.

BRIEF DESCRIPTION OF THE DRAWINGS

Figure 1 is a block diagram of a packet switched system constructed in accordance with an embodiment of the present invention.

Figure 2 is a block diagram of a multiport switch constructed in accordance with an embodiment of the present invention and used in the packet switched system of Figure 1.

Figure 3 is a schematic depiction of a switch subsystem of the multiport switch of Figure 3, constructed in accordance with an embodiment of the present invention.

Figure 4 is a block diagram of a single output queue of the switch subsystem of Figure 4, constructed in accordance with an embodiment of the present invention.

Figure 5 is a detail of a first type of output queue in accordance with an embodiment of the present invention.

Figure 6 is a detail of a second type of output queue in accordance with an embodiment of the present invention.

Figure 7 is a detail of an overflow area of the external memory, configured in accordance with an embodiment of the present invention.

Figure 8 is a block diagram of a linked list data structure employed in the present invention.

Figure 9 schematically depicts a frame buffer header format in accordance with an embodiment of the present invention.

Figure 10 is a detail of the multicopy, reclaim and free buffer pool area of the switch subsystem of Figure 4, constructed in accordance with an embodiment of the present invention.

Figure 11 is block diagram of a free buffer pool structure constructed in accordance with an embodiment of the present invention.

Figure 12 is a block diagram of a multicopy queue configured in accordance with an embodiment of the present invention.

Figure 13 is a schematic representation of the multicopy cache constructed in accordance with an embodiment of the present invention.

Figure 14 is a block diagram of a queuing block of the buffer manager of the switch subsystem and a port vector FIFO, constructed in accordance with an embodiment of the present invention.

Figure 15 is a block diagram of an output queue constructed in accordance with another embodiment of the present invention.

DETAILED DESCRIPTION OF ILLUSTRATIVE EMBODIMENTS

The present invention will be described with the example of a switch in a packet switched network, such as an Ethernet (IEEE 802.3) network. It will become apparent, however, that the present invention is also applicable to other packet switched systems, as described in detail below, as well as to other types of systems in general.

Figure 1 is a block diagram of an exemplary system in which the present invention may be advantageously employed. The exemplary system 10 is a packet switched network, such as an Ethernet network. The packet switched network includes an integrated multiport switch (IMS) 12 that enables communication of data packets between network stations. The network may include network stations having different configurations, for example twenty-four (24) 10 megabit per second (Mbps) network stations 14 that send and receive data at a network data rate of 10 Mbps, and two 100 Mbps network stations 22 that send and receive data packets at a network speed of 100 Mbps. Hence, the switch 12 selectively forwards data packets received from the network stations 14 or 22 to the appropriate destination based upon Ethernet protocol.

According to the disclosed embodiment, the 10 Mbps network stations 14 send and receive data packets to and from the switch 12 via a media 17 and according to half-duplex Ethernet protocol. The Ethernet protocol ISO/IEC 8802-3 (ANSI/IEEE Std. 802.3, 1993 Ed.) defines a half-duplex media access mechanism that permits all stations 14 to access the network channel with equality. Traffic in a half-duplex environment is not distinguished or prioritized over the medium 17. Rather, each station 14 includes an Ethernet interface card that uses carrier-sense multiple access with collision detection (CSMA/CD) to listen for traffic on the media. The absence of network traffic is detected by sensing a deassertion of a receive carrier on the media. Any station 14 having data to send will attempt to access the channel by waiting a predetermined time after the deassertion of a receive carrier on the media, known as the interpacket gap interval (IPG). If a plurality of stations 14 have data to send on the network, each of the stations will attempt to transmit in response to the sensed deassertion of the receive carrier on the media and after the IPG interval, resulting in a collision. Hence, the transmitting station will monitor the media to determine if there has been a collision due to another station sending data at the same time. If a collision is detected, both stations stop, wait a random amount of time, and retry transmission.

The 100 Mbps network stations 22 preferably operate in full-duplex mode according to the proposed Ethernet standard IEEE 802.3x Full-Duplex with Flow Control – Working Draft (0.3). The full-duplex environment provides a two-way, point-to-point communication link between each 100 Mbps network station 22 and the switch 12, where the switch 12 and the respective stations 22 can simultaneously transmit and receive data packets without collisions. The 100 Mbps network stations 22 each are coupled to network media 17 via 100 Mbps physical (PHY) devices 20 of type 100 Base-TX, 100 Base-T4, or 100 Base-FX. The switch 12

includes a media independent interface (MII) 24 that provides a connection to the physical devices 20. The 100 Mbps network stations 22 may be implemented as servers or routers for connection to other networks.

As shown in Figure 1, the network 10 includes a series of switch transceivers 26 that perform time division multiplexing and time division demultiplexing for data packets transmitted between the switch 12 and the 10 Mbps stations 14. A magnetic transformer module 19 maintains the signal waveform shapes on the media 17. The switch 12 includes a transceiver interface 18 that transmits and receives data packets to and from each switch transceiver 16 using a time-division multiplexed protocol across a single serial non-return to zero (NRZ) interface 23. The switch transceiver 16 receives packets from the serial NRZ interface 23, demultiplexes the received packets, and outputs the packets to the appropriate end station 14 via the network media 17. According to the disclosed embodiment, each switch transceiver 16 has four independent 10 Mbps twisted-pair ports and uses 4:1 multiplexing across the serial NRZ interface enabling a four-fold reduction in the number of PINs required by the switch 12.

The switch 12 contains a decision making engine, switching engine, buffer memory interface, configuration/control/status registers, management counters, and MAC (media access control) protocol interface to support the routing of data packets between the Ethernet ports serving the network stations 14 and 22. The switch 12 also includes enhanced functionality to make intelligent switching decisions, and to provide statistical network information in the form of management information base (MIB) objects to an external management entity, described below. The switch 12 also includes interfaces to enable external storage of packet data and switching logic in order to minimize the chip size of the switch 12. For example, the switch 12 includes a synchronous dynamic RAM (SDRAM) interface 34 that provides access to an external memory 36 for storage of received frame data, memory structures, and MIB counter information. The memory 36 may be an 80, 100 or 120 MHz synchronous DRAM having a memory size of 2 or 4 Mb.

The switch 12 also includes a management port 30 that enables an external management entity to control overall operations of the switch 12 by a management MAC interface 32. The switch 12 also includes a PCI interface 26 enabling access by the management entity via a PCI host and bridge 28. Alternatively, the PCI host and bridge 28 may serve as an expansion bus for a plurality of switch devices 12.

The switch 12 includes an internal decision making engine (Figure 2) that selectively transmits data packets received from one source to at least one destination station. The internal decision making engine may be substituted with an external rules checker. The switch 12 includes an external rules checker interface (ERCI) 40 that allows use of an external rules checker 42 to make frame forwarding decisions in place of the internal decision making engine. Hence, frame forwarding decisions can be made either by the internal switching engine or the external rules checker 42.

The switch 12 also includes an LED interface 44 that clocks out the status of conditions per port and drives LED external logic 46. The LED external logic 46, in turn, drives LED display elements 48 that are human readable. An oscillator 38 provides a 40 MHz clock input for the system functions of the switch 12.

Figure 2 is a block diagram of the integrated multiport switch (IMS) 12 of Figure 1. The switch 12 includes twenty-four (24) 10 Mbps media access control (MAC) ports 50 for sending and receiving data packets in half-duplex between the respective 10 Mbps network stations 14 (ports 1-24), and two 100 Mbps MAC ports 53 for sending and receiving data packets in full-duplex between the respective 100 Mbps network stations

(ports 25, 26). As described above, the management interface 30 also operates according to MAC layer protocol (port 0). Each of the MAC ports 50, 53 and 30 has a receive first in-first out (FIFO) buffer 52 and transmit FIFO 54. Data packets from a network station are received by the corresponding MAC port and stored in the corresponding receive FIFO 52. The received data packet is output from the corresponding receive FIFO 52 to the external memory interface 34 for storage in the external memory 36.

The header of the received packet is also forwarded to a decision making engine, either internal rules checker 58 or the external rules checker interface 40, to determine which MAC ports will output the data packet. Specifically, the packet header is forwarded to the internal rules checker 58 or the external rules checker interface 40, depending on whether the switch 12 is configured to operate using the internal rules checker 58 or the external rules checker 42. The internal rules checker 58 and external rules checker 42 provide the decision making logic for determining the destination MAC port for a given data packet. The decision making engine may thus output a given data packet to either a single port, multiple ports, or all ports (i.e., broadcast). For example, each data packet includes a header having source and destination address, where the decision making engine may identify the appropriate output MAC port based upon the destination address. Alternatively, the destination address may correspond to a virtual address that the appropriate decision making engine identifies as corresponding to a plurality of network stations. Alternatively, the received data packet may include a VLAN (virtual LAN) tagged frame according to IEEE 802.1d protocol that specifies another network (via a router at one of the 100 Mbps stations 22) or a prescribed group of stations. Hence, either the internal rules checker 58 or the external rules checker 42 via the interface 40 will decide whether a frame temporarily stored in the buffer memory 36 should be output to a single MAC port or multiple MAC ports.

Use of the external rules checker 42 provides advantages such as increased capacity, a random-based ordering in the decision queue that enables frame forwarding decisions to be made before the frame is completely buffered to external memory, and enables decisions to be made in an order independent from the order in which the frames were received by the switch 12.

The decision making engine (i.e., internal rules checker 58 or the external rules checker 42) outputs a forwarding decision to a switch subsystem 56 in the form of a port vector identifying each MAC port that should receive the data packet. The port vector from the rules checker includes the address location storing the data packet in the external memory 36, and the identification of the MAC ports to receive the data packet for transmission (e.g., MAC ports 0-26). The switch subsystem 56 fetches the data packet identified in the port vector from the external memory 36 via the external memory interface 34, and supplies the retrieved data packet to the appropriate transmit FIFO 54 of the identified ports.

Additional interfaces provide management and control information. For example, a management data interface 59 enables the switch 12 to exchange control and status information with the switch transceivers 16 and the 100 Mbps physical devices 20 according to the MII management specification (IEEE 802.3u). For example, the management data interface 59 outputs a management data clock (MDC) providing a timing reference on the bidirectional management data IO (MDIO) signal path.

The PCI interface 26 is a 32-bit PCI revision 2.1 compliant slave interface for access by the PCI host processor 28 to internal IMS status and configuration registers 60, and access external memory 36. The PCI interface 26 can also serve as an expansion bus for multiple switch devices. The management port 30 interfaces

to an external MAC engine through a standard seven-wire inverted serial GPSI interface, enabling a host controller access to the switch 12 via a standard MAC layer protocol.

Figure 3 depicts the switch subsystem 56 of Figure 2 in more detail according to an exemplary embodiment of the present invention. Other elements of the multiport switch 12 of Figure 2 are reproduced in Figure 3 to illustrate the connections of the switch subsystem 56 to these other elements. The switch subsystem 56 contains the core switching engine for receiving and forwarding frames. The main functional blocks used to implement the switching engine include: a port vector FIFO 70, a buffer manager 72, a plurality of port output queues 74, a management port output queue 75, an expansion bus port output queue 77, a free buffer pool 104, a multicopy queue 90, a multicopy cache 96 and a reclaim queue 98. The operation and structure of these functional blocks will be described in more detail, but a brief overview of the switch subsystem 56 of Figure 3 is first presented to provide context to the later discussion of the individual elements.

There are two basic types of frames that enter the multiport switch 12 from the ports: unicast frames and multicopy frames. A unicast frame is a frame that is received at a port which is to be transmitted by the multiport switch 12 to only one other port. By contrast, a multicopy frame is a frame that is received at one port for transmission to more than one port. In Figure 3, each port is represented by a separate MAC 50, having its own receive FIFO 52 and transmit FIFO 54.

Frames, whether unicast or multicopy, are received by the internal MAC engines 50. When the frame packet is received at the port, it is placed in the receive FIFO 52. Each frame has a header, which is provided to a rules checker, either the internal rules checker 58 or the external rules checker 42. The rules checker 42 or 58, based on the information in the header, determines from where the frame packet will be cast, i.e., through which port or ports will the frame packet be transmitted.

At the same time as the rules checker 42 or 58 is making its forwarding determination, the buffer manager 72 obtains a free buffer pointer from the free buffer pool 104. This free buffer pointer is the location in external memory 36 at which the frame will be stored by the receive FIFO 52. Once the free buffer pointer is obtained from the free buffer pool 104 by the buffer manager 72, the buffer pointed to by the free buffer pointer is no longer considered free. The frame data is transferred over data bus 80 from the receive FIFO 52 to the external memory 36 in a direct memory access (DMA) transaction. The frame is stored in the location pointed to by the free buffer pointer obtained from the free buffer pool 104, although a number of other buffers may be used to store a frame, as will be described.

In addition to the header data, the rules checker 42 or 58 also receives the free buffer pointer from the buffer manager 72. This free buffer pointer is now referred to as a frame pointer since it points to the memory location in the external memory 36 where the frame is stored. The rules checker 42 or 58 uses the header information to make the forwarding decision and generate a forwarding instruction in the form of a "port vector". In the exemplary illustrated embodiment, the port vector is a 28-bit vector with a bit set for each output port to which the frame should be forwarded. Assume for this overview example that the received frame is a unicast frame. Accordingly, only one bit is set in the port vector generated by the rules checker 42 or 58. The bit that is set in the port vector corresponds to a particular one of the ports.

The rules checker 42 or 58 places the port vector and the frame pointer (as well as a control opcode and a VLAN index) into the port vector FIFO 70. The port vector is examined by the port vector FIFO 70 to

determines into which particular output queue 74 (or queues) the frame pointer associated with the port vector should be input. The port vector FIFO 70 places the frame pointer into the top of the appropriate output queue 74. This queues the transmission of the frame.

At some point in time, the frame pointer reaches the bottom of the output queue 74 after passing through the output queue 74. The buffer manager 72 takes the frame pointer when it arrives at the bottom of the output queue 74 and passes the frame pointer to the appropriate transmit FIFO 54 of the correct port via frame pointer read bus 86. This schedules the transmission of the frame. The frame data is read in a DMA transaction from the location in external memory 36 pointed to by the frame pointer, is placed in the appropriate transmit FIFO 54 and then transmitted.

A multicopy transmission is similar to the unicopy transmission, except that the port vector has multiple bits set, designating the multiple ports from which the frame will be transmitted. The frame pointer is placed into each of the appropriate output queues 74 and transmitted from the corresponding transmit FIFOs 54.

The buffer manager 72 uses the special control queues, i.e., the free buffer pool 104, the multicopy queue 90, and the reclaim queue 98, and the multicopy cache 96 to manage the process of allocating buffers to store received frames and retrieving buffers for re-use once the frame has been transmitted to its designated output port(s). The buffer manager 72 also maintains "overflow" regions in external memory 36 for the output queues 74 and the control queues 104, 90 and 98, as will be described in more detail later.

With this operational overview serving as background, the individual sections and various aspects of the switch subsystem 56 will now be discussed in more detail. The first of these aspects that will be described is the structure of the various output queues 74 of the present invention. In addition to the output queues 74 designated for the 10 Mb/s output ports and the 100 Mb/s output ports, an output queue 75 is provided for the management port 30, and an output queue 77 is provided for the expansion port 26. These output queues 75, 77 have the same external structure as the output queues 74, but different internal configurations, as will be described.

Figure 4 is a block diagram of the external structure of an output queue 74 in accordance with an embodiment of the present invention. As is apparent from Figure 4, the output queue 74 of the present invention has a three-part configuration. For highest performance it is preferable to keep all of the queuing structure on the chip (referring to the multiport switch 12), but the real estate on a chip is very expensive. This presents a dilemma when the chip is designed to switch, and needs to queue, a large number of entries. The present invention solves this dilemma by providing a single output queue that includes a high performance, low capacity section that is on-chip, and an overflow area that is off-chip. The overflow area allows the queue to serve as a large capacity queue as needed, albeit with a relatively lower performance than the on-chip portion.

A single logical output queue 74 of the present invention, according to the embodiment of Figure 4, has three physical sections. These include an output queue write side 76, an output queue read side 78, and an output queue overflow area (generally designated as 110) located in the external memory 36. Access to the external memory 36 for all of the output queues 74 is through the external memory interface 34, as described earlier. The present invention takes advantage of the bursting nature of current external memories, so that the data (e.g., frame pointers) is sent on and off the chip to the overflow queue area 110 in bursts over the bus 84 connecting the chip 12 to the external memory 36.

The output queue write side 76 and the output queue read side 78 are located on the chip 12. The write side 76 and the read side 78 are considered to be small, expensive resources. By contrast, the overflow area 110, forming the third part of the output queue 74, is large and inexpensive. The write side 76 and the read side 78 provide high performance, while the path through the overflow area provides a low-performance, large capacity path.

In operation, the output queue write side 76 receives an entry. In the exemplary embodiment of a multiport switch 12 according to the present invention, the entry is a frame pointer that points to the first buffer in external memory in which the first 256 bytes of a frame are stored. It should be apparent to those of skill in the art, however, that the output queue structure 74 is not limited to frame pointers as entries, but is widely applicable to queue other types of entries, both in multiport switches and in other technologies.

After the entry flows through and reaches the bottom of the output queue write side 76, control logic associated with the output queue 74 makes a decision as to what to do with the entry. If there is space in the output queue read side 78, and the overflow area 110 for that output queue 74 is empty, then one or more entries are passed directly from the output queue write side 76 to the output queue read side. This passing of the entry or entries directly from the write side 76 to the read side 78 is performed entirely on the chip 12, and is therefore a low-latency, fast flow-through of an entry.

If the output queue read side 78 is full, and there is at least a burst-size amount of data (e.g., 16 bytes worth of entries) in the output queue write side 76, then the data is written in a burst fashion into the overflow area 110 for that output queue 74. If the output queue read side 78 is full, but there is not yet a burst-size amount of data in the output queue write side 76, then the entry remains in the output queue write side and nothing further is done. Eventually, the output queue read side 78 will empty, and when the output queue read side 78 has enough space to accommodate a burst-size amount of data, and there is data in the overflow area 110, a burst of data is provided from the overflow area 110 into the output queue read side 78.

In the output queue structure, the read side 78 is acting most like a traditional queue, because it is from this portion that entries are taken, one by one. The output queue write side 76 mostly serves a collection function to assemble the data into bursts for writing to the external memory 36. Hence, the present invention transforms single events (placing an entry into the output queue 74) into a burst event. The write side 76 allows an accumulation of data to then burst, if necessary, to the overflow area 110 in the external memory 36. The overflow area 110 provides inexpensive storage in times of congestion, rather than dedicating expensive chip resources to a function that will only be needed on relatively rare occasions. Even though the present invention utilizes an overflow area 110 that is off-chip, the accessing of this area 110 is performed in a manner that is efficient, by bursting a number of bytes of information at a time. This is in contrast to conventional queuing structures in which single entries are written and read to and from the queue.

During operation, if there is a lot of entries arriving at the output queue 74, these entries are placed into the overflow area 110 to avoid overflowing the on-chip queue 78. Hence, the discarding of frames is largely prevented with the queue structure of the present invention. Also, the total amount of memory dedicated to the overflow areas 110 may be readily changed by changing the size of the external memory 36. Furthermore, the sizes of the individual specific overflow areas 110 are programmable to customize the queue sizes, without impacting the performance of the output queues 74.

Typically, a queue is an ordered structure with a first-in, first-out arrangement. In some types of queues, however, such as the reclaim queue 98 and the free buffer pool 104, the order of entries does not matter. If it is possible to send data directly from the write side 100 to the read side 102, the present invention permits information to be directly sent this route, bypassing the overflow area for the queue. This is permitted even if there is information in the associated overflow area, as long as the information is not order-sensitive. For example, the reclamation of buffers is not order-sensitive, since any order in which the buffers are eventually returned to the free list in the free buffer pool 104 after the buffer is no longer needed to store the frame is acceptable. Hence, in order to avoid incurring the bandwidth of a write to the overflow area 110 for the reclaim queue 98 in the external memory 36 when the data is not order-sensitive, the information is passed directly from the write side 100 to the read side 102, assuming the read side 102 has room for more entries. The reclaim queue 98 is an example of a type of queue that queues data which is not order-sensitive. However, there are many other types of data in different applications that are also not order-sensitive, so that this feature of the present invention finds utility in queues that queue these other types of data.

In the multiport switch of an exemplary embodiment of the present invention as depicted in Figures 1 and 2, there are twenty-eight output queues (each associated with an output port): twenty-four for the 10 Mb/s user ports, two for the 100 Mb/s server ports, one for the management port and one for the expansion bus port. The output queues 74, 75 and 77 provide temporary storage for frame pointers when they are queued for transmission. Queuing takes the form of the port vector FIFO 70 writing frame pointers into the various output queues 74, 75 and 77 indicated in a forwarding port vector.

In certain preferred embodiments of the invention, the various output queues 74, 75 and 77 contain several or all of the following fields: unicopy bit, frame pointer, control opcode or control signals, and VLAN (virtual local area network) index. The unicopy bit flags a frame which is to be forwarded to only one output port. The frame pointer points to the frame in external memory 36. The control opcode identifies special information about the frame (i.e., newly learned frame, etc.). The control signals use information from the control opcode to indicate how the ports will handle frames before transmission. The VLAN index provides the reference to a VLAN tag which should be inserted (if necessary) into the outgoing frame. However, these fields are exemplary only, as the present invention is applicable to other output queues with different types of fields.

The internal structure of an exemplary embodiment of a first type of output queue 74, the 10 Mb/s port output queue, is depicted in Figure 5. The 10 Mb/s output queues 74 hold entries for frames to be forwarded to the 10 Mb/s ports. The output queue write sides 76 for these queues hold thirty-two entries and the output queue read sides 78 hold sixteen entries in the exemplary illustrated embodiment, although other sizes are contemplated and within the scope of the invention. Each entry in a 10 Mb/s output queue 74 comprises a unicopy bit and a frame pointer (14 bits). In the exemplary embodiment of the multiport switch of the present invention, the VLAN index is not necessary because there is no VLAN tagging on 10 Mb/s ports.

The internal structure of an exemplary embodiment of a second type of output queue 74, the 100 Mb/s port output queue, is depicted in Figure 6. The 100 Mb/s port output queues hold entries for frames to be forwarded to the 100 Mb/s ports. The output queue write side 76 holds sixty-four entries in this type of output queue, and the output queue read side holds sixteen entries. Each entry comprises a VLAN index, a partial control opcode (bits 4-0), a unicopy bit and a frame pointer.

An alternative embodiment of the output queue structure is depicted in Figure 15 in block diagram form. The output queue 190 may be substituted for the output queues 74, 75 and 77 in the switch 12 of Figure 3. The output queue 190 according to the embodiment of Figure 15, provides prioritization of frames in a space saving, and therefore, cost-effective manner.

The queue 190, like the queues 74, 74 and 77, is a split queue with a queue write side 192. However, the queue 190 has a plurality of read sides 194a - n, each of which will queue entries from the write side 192. In the embodiment of Figure 15, each of the read sides 194a-n is assigned a different priority level, and will therefore handle different priority level traffic. For example, read side 194a may be assigned the highest priority level of the read sides, and will then queue those entries exiting the write side 192 that have a matching priority level, i.e., the highest priority level. Similarly, queue read side 194n may be assigned the lowest priority level, and will queue entries that have a matching low priority level. In certain embodiments, the higher priority queue read sides 194 are smaller in queuing capacity than the lower priority queue read sides 194, since the high priority read sides should be kept small to provide a low latency for the high priority entries.

When the exemplary embodiment of the invention is employed in the network switch 12 of Figure 3, the entries are frame pointers that point to a location in the external memory 36 in which the frame is stored. The frame pointers are then provided with additional bits, for example three bits, that encode the priority of the frame pointed to by the frame pointer. This priority may be encoded, for example, by the rules checker 42 or 58 when making the forwarding decision. The priority bits accompany the frame pointer as it is passed around the switch 12. The priority of the frame may be based on the time sensitivity of the data, or the content of the frame, for instance, although it may be based on other factors as well.

The passing of the frame pointer to the appropriate priority level read side 194 is controlled by control logic 196, which examines the priority bits of each entry exiting from the queue write side 192. Upon determining the priority of an entry, the control logic 196 checks the available capacity of the queue read side 194 with the priority level that matches the priority level of that entry. If there is available capacity, the control logic 196 loads the entry into that queue read side 194. In certain embodiments, if there is no available capacity in the matching queue read side 194, then the entry is placed into the overflow area 110 in the external memory 36 through the external memory interface 34. When that particular queue read side 194 regains capacity, then the entry is retrieved from the overflow area 110 and placed into the queue read side 194.

In certain embodiments of the invention, high priority entries are discarded if the high priority queue read side (194a, for example) does not have available capacity. This discarding of the entry (and thereby the frame copy) instead of placing the entry into the overflow area 110 for that queue 190, recognizes that high priority traffic that is delayed by the amount of time it takes to write the entry to the overflow area and retrieve the entry therefrom, is unlikely to be useful as high priority data. Rather than send this frame, which is most likely stale, these embodiments of the invention discard the frame and thereby seek to alleviate the high priority congestion. These embodiments of the invention are not limited to those that discard only the highest priority frames. A threshold priority level may be provided to the control logic 196 so that frames of a priority level above the threshold level will be discarded if the matching queue read sides 194 have no available capacity.

In certain embodiments of the invention the control logic 196, upon determining that the queue read side 194 with the matching priority level does not have any available capacity, attempts to load the entry into the

next lower priority level queue read side 194. This method attempts to avoid sending the entry to the overflow area 110 and still try to queue the entry with a priority as close as possible to the actual priority of the entry.

The control logic 196 is readily implemented by one of ordinary skill in the art given the descriptions of the functionality provided above.

The MAC 50 polls the queue read sides 194a-n and obtains queue entries from the highest priority queue read sides 194 before beginning to empty the lower priority queue read sides 194. For example, assume that queue read side 194a is the highest priority queue read side, queue read side 194b the next highest, and queue read side 194n is the lowest priority queue read side. The MAC 50 will always service and empty the queue read side 194a prior to servicing queue read side 194b. As long as there is an entry in queue read side 194a, the MAC 50 will not begin emptying any other queue read side. If queue read side 194a is empty, and queue read side 194b has an entry, the MAC 50 will service this queue read side 194 b before servicing any other lower priority queue read sides 194c-n until the queue read side 194b is completely empty, and so on.

The priority queue structure according to the embodiment of Figure 15 has a reduced size in comparison to arrangements in which the entire queue is replicated to produce queues of different priorities. This greatly adds to the size and cost of a device, especially a multiport device. Although the queue structure of Figure 15 finds great utility in its application to a multiport switch, such as switch 12 of Figure 3, those of ordinary skill in the art will recognize that the structure is applicable in a wide range of other types of devices in which data is to be queued in a prioritized manner. For example, in computing devices that have on-board queues, certain data may be considered high priority data. Saving space dedicated to the on-board queues, and still provide prioritization of the data, is a desirable feature. Another example are phone devices that process both voice and computer traffic. The voice traffic may be deemed higher priority traffic, and the computer traffic lower priority.

An exemplary map of the external memory 36 is depicted in Figure 7. The overall capacity of the external memory 36 may be, for example, 4 Mb, although other capacity memories are employed in different embodiments. The use of an external memory 36 for the overflow areas according to the present invention permits increasing or decreasing the size of the output queues by simply changing the external memory. This is an advantage over systems in which the queue structure is entirely on the chip, as the overall queuing capacity is set at manufacture of the chip.

To satisfy the storage requirements of the switch 12; an exemplary embodiment of the external memory 36 allocates space for the following areas: free buffer pool overflow 120, reclaim queue overflow 122, multicopy queue overflow 124, management port output queue overflow 126, individual output queue overflows 128 for each of the 10 Mb/s and 100 Mb/s destination ports, expansion bus port output queue overflow 130, the MIB counters 132, and the global frame buffer pool 134.

The BASE Address for the entire memory region is programmable in a memory base address register among the registers 60 on the chip. The BASE Address for each area in the external memory map is programmable in the register set. No length register is required, the length for a given area being equal to the area from that area's BASE Address to the BASE Address of the next area in the mapping.

Since the length (and therefore capacity) of each of the individual overflow areas is programmable, the overall capacity of each queue is programmable. This feature of the present invention permits customization of the switch to provide particular output queues with increased capacity, as needed.

The following overflow areas store entries that do not fit into the control queues on the chip 12 are therefore placed into the external memory 36. The free buffer pool overflow area 120 stores the address pointers to currently unused buffers in the global frame buffer pool 134. The reclaim queue overflow area 122 stores frame pointers to linked-list chains that are no longer needed. The multicopy queue overflow area 124 stores frame pointers with copy numbers ">=1" (for queued frame pointers) and frame pointers with copy numbers "-1" (for successfully transmitted frames)

The following overflow areas store entries for output queues which do not fit on-chip. The management port output queue overflow area 126 stores frame pointers awaiting transmission to the management port. Output queue overflow areas 128 store frame pointers awaiting transmission to the appropriate 10 Mb/s port or 100 Mb/s port. The expansion bus port output queue overflow area 130 stores frame pointers awaiting transmission to the expansion bus port.

The MIB counter region 132 contains all the per port statistics which are updated periodically by the switch 12. The switch 12 maintains 8-bit and 16-bit counters on-chip for storing MIB statistics. The switch 12 updates the 32-bit or 64-bit MIB counters in external memory 36 with the frequency required to prevent loss of MIB data.

The global frame buffer pool 134 contains buffers in linked-lists which store received frame data. At any given time, these linked lists contain valid frame data, obsolete buffers which will be returned by the buffer manager 72 to the free buffer pool 104 or are owned by the PCI host processor 28.

Referring now to Figure 8, frame data received from any MAC port or the PCI bus is stored in external memory 36 in a linked-list data structure format in an exemplary embodiment of the present invention. The buffers 140 used to create the linked-list are 256 bytes in length, although other sized buffer lengths are employed in different embodiments of the invention. Address pointers to each of these buffers 140 are stored by the free buffer pool 104 in the switch 12.

As a frame is received at one of the ports of the switch 12, the buffer manager 72 requests address pointers from the free buffer pool 104 for linking buffers 140 to store the frame. The address pointer to the first buffer in external memory 36 that stores the frame becomes the frame pointer for that frame. The frame pointer is used in the switch subsystem 56 for queuing frames to be transmitted.

The buffers 140 are chained together by address pointers in each buffer header 142 that indicate the location of the next buffer in memory. The buffer headers 142 also contain other information about the frame data contained in the buffer 140. The first buffer's header is 12 bytes, as depicted in the exemplary buffer header format of Figure 9a. Each subsequent buffer's header is 4 bytes, as depicted in Figure 9b. The external memory bursts are 2 banks x 16-bytes long, so actual frame storage capacity in each buffer is $256B - 16B = 240B$.

As depicted in Figures 9a and 9b, the first and subsequent buffer header formats contain the following fields:

Buffer Format Bit: indicates what buffer format is in use. A one indicates the first buffer's format, which is 12 bytes in length. A zero indicates a subsequent buffer's format, which is 4 bytes. It is used for each of the remaining buffers when chaining buffers.

E Bit (End of Frame Marker): indicates this is the last buffer for a frame. When the E bit is set, there are no more buffers in the chain.

C Bit (CRC Error Detected): indicates a CRC error was detected by the receiver. When the C Bit is detected, the transmit function will purposely transmit an inverted CRC.

L Bit (Alignment Error): indicates a Frame Alignment Error was detected (along with a CRC Error) in the receive frame.

O Bit (Receive FIFO Overflow): indicates the receive FIFO overflowed and the data in the buffer may not be valid.

Buffer Length: the total number of bytes which are valid in the data field of the buffer beginning with the first byte after the buffer header. This length should not include the Offset Byte value.

Next Buffer Pointer: the pointer to the next buffer. The next buffer pointer is not valid when the E Bit is set.

Offset Byte Count: indicates where the first byte of the frame starts in the frame data section of the buffer. An offset of zero means the data will begin at the first byte after the buffer header 142. An offset of zero indicates frame data will begin at the byte following the 16th byte in the buffer. For non-zero values of offset, frame data will begin following 16B + Offset from the beginning of the buffer. The transmit function will skip over the number of bytes indicated in the offset Byte field.

P Bit (Port Type): indicates the port type of the incoming receive frame. A zero indicates a 10 Mb/s port and a one indicates a 100 Mb/s port. This bit is used by the host 28 in conjunction with the time stamp field when it programs the switch 12 to forward frames to the expansion bus before the frame is completely received and buffered to external memory 36.

T Bit: indicates the received frame type: tagged or untagged. A one indicates a tagged frame and the VLAN Identifier field contains the received VLAN ID. A zero indicates an untagged frame and the VLAN ID is not valid.

Receive Port Number: the number of the port from which the frame was received.

VLAN Identifier: the VLAN ID received from a "tagged" port. If the frame was received from an untagged port, this field is invalid.

R Bit (Recalculate CRC): indicates the CRC needs to be stripped and recalculated at the transmit function. The switch 12 sets this bit when a tagged frame is received. In addition, if the host 28 modifies a frame's contents, the host 28 should set this bit. When the switch 12 transmits a frame, it will examine this bit to determine whether to transmit the existing CRC or strip and recalculate the CRC.

A Bit (Append CRC): indicates that there is no CRC at the end of the frame data. The host can create a frame in memory (without a CRC) then set this bit. The switch 12 will generate and append a CRC when transmitting the frame. If the A Bit is set, the frame length should not include CRC.

F Bit (Format Bit): identifies the Frame Length/Time Stamp field. A zero indicates the field is the time stamp of the incoming frame. A one indicates the field is the frame length of the received frame.

Frame length/time stamp: dependent on F Bit. If F Bit is cleared, this field represents the time stamp from the beginning of the received frame. The time stamp has a resolution of 1 μ s. If the F Bit is set, indicates the total length of the received frame including CRC and any received VLAN Tag. When a frame is received, the switch 12 marks this field with the time stamp (from the timer register). If the host 28 has programmed the switch 12 to forward expansion bus frames before the frame has been completely received, it can use the time stamp (along with the speed of the receive port) to gauge how much data it can fetch from external memory 36 without over-reading the frame data. Once the entire frame has been received, the switch 12 writes the frame length into this field and sets the F Bit.

Copy number: used to indicate the number of copies successfully queued for transmission by the port vector FIFO 70. This field is used to store the copy number for a frame pointer if the buffer manager 72 needs to make space in the multicopy cache 96 for new entries.

Figure 10 is a detailed depiction of some of the elements of the switch subsystem 56 of Figure 3. These elements are used to provide the buffers for storage of frames, and to reclaim these buffers and make them available for use again once the buffers are no longer needed for storage of the frame. As described earlier, each output queue 74, 75 (except output queue 77) passes frame pointers to the buffer manager 72, which schedules transmission of the frames pointed to by the frame pointers. The buffer manager 72 controls the following functions: 1) managing the internal busses of the switch 12; 2) facilitating queuing/dequeuing frame pointers to/from the output queues 74; 3) managing the control queues 90, 98 in order to locate and return buffers to the free buffer pool 104; 4) controlling the flow of data to and from the external memory 36; and 5) maintaining the memory structures, including MIBs and overflow areas. The buffer manager 72 contains a scheduler function for allocating all accesses to external memory 36. These accesses include 1) writing received frame data to memory buffers 140, 2) reading frame data from memory buffers 140 for transmission and 3) maintaining (i.e., writing and reading) the frame pointers in each of the overflow areas for the output queues 74 and control queues 90, 98 and 4) updating MIB counters.

After the buffer manager 72 has copied a given frame pointer to all the appropriate output queue(s) 74, 75, the port vector FIFO 70 calculates the number of copies (the "copy number") and places the frame pointer and the copy number into the write side 92 of the multicopy queue 90. The copy number may be "0", indicating that the frame should not be forwarded, a "1", indicating a unicity transmission, or a number ">1", indicating a multicopy transmission. These three cases are described below.

When the copy number is "0", which means that the frame pointer has a null forwarding port vector with no bits set, the port vector FIFO 70 passes the frame pointer directly to the write side 100 of the reclaim queue 98. When the buffer manager 72 services the reclaim queue 98, as will be described, the buffer manager 72 breaks down the linked-list chain of buffers and returns the address pointer for each "free" buffer to the write side 106 of the free buffer pool 104.

When the copy number is "1", a unicity transmission, the port vector FIFO 70 copies the frame pointer, control signals/control opcode and the VLAN index to the output queue 74 of the appropriate port. The port vector FIFO 70 sets the unicity bit in the output queue 74 (see Figures 5 and 6) to indicate that this is a single transmission. When the buffer manager 72 reads the frame pointer and the unicity bit from the output queue 74 of the port, it schedules the transmission as discussed previously. The buffer manager 72 uses the

frame pointer to locate the first buffer in external memory 36 in which the frame is stored. The buffer manager 72 reads the buffer header from this first buffer, captures data from the first buffer and places this data in the appropriate MAC transmit FIFO 54. The links to subsequent buffers, assuming that the frame spans multiple buffers, provides the buffer manager 72 with the address to find and transmit all of the buffers in the chain for that frame. Once the data has been placed in the FIFO 54 for transmission, the buffer becomes obsolete and is returned to the free buffer pool 104 for eventual re-assignment to store data of another frame.

When the copy number is greater than 1, the port vector FIFO 70 copies the frame pointer, VLAN index and control signals/control opcode to each of the appropriate output queues 74. (When referring to queues 74, reference is also made to queues 75, 77). The port vector FIFO 70 clears the unicopy bit for the appropriate frame pointers in the output queues 74 and places the frame pointer with a copy number of ">1" into the write side 92 of the multicopy queue 90.

Whenever the buffer manager 72 reads a frame pointer and a cleared unicopy bit from one of the output queues 74, the buffer manager 72 schedules the transmission of the frame, but also checks the multicopy cache 96 for an entry with a frame pointer having a copy number of "1". If a frame pointer with a copy number of "1" is found in the multicopy cache 96, then the buffer manager 72 schedules the frame for transmission and reclaims the buffers during transmission in the same manner as in the unicopy transmission of a frame. However, if the frame pointer is not in the multicopy cache 96 or the copy number of the frame pointer in the multicopy cache 96 is greater than 1, then the buffer manager 72 transmits the frame but does not reclaim the buffers. After successful transmission, the buffer manager 72 places a copy of the frame pointer, along with a copy number of "-1" into the write side 92 of the multicopy queue 90.

Each time a multicopy frame is transmitted, the buffer manager 72 places a copy of the frame pointer into the multicopy queue 90, provided the buffer manager 72 did not find the frame pointer in the multicopy cache 96 with a copy number of "1". Hence, at any given time, the multicopy queue 90 may contain the frame pointer with a copy number that is > "1" and/or several copies of the same frame pointer, each with a copy number of -1.

The buffer manager 72 constantly services the multicopy queue 90 and the multicopy cache 96 in order to reclaim obsolete buffers. When it services the multicopy queue 90 and reads a frame pointer with a copy number ">1", the buffer manager 72 attempts to place this new entry (frame pointer and copy number) into the multicopy cache 96. If the multicopy cache 96 is full, the buffer manager 72 makes space for the new frame pointer. The buffer manager 72 reads an "older" multicopy cache entry, updates the copy number for this entry in its buffer header in external memory 36, then clears the entry from the multicopy cache 96. As room becomes available in the multicopy cache 96, the buffer manager 72 is able to place the new entry from the multicopy queue 90 into the multicopy cache 96.

When the buffer manager 72 services the multicopy queue 90 and reads a frame pointer with a copy number of "-1", it searches the multicopy cache 96, looking for a matching frame pointer address with a copy number ">=1" to decrement or delete. If the buffer manager 72 finds a frame pointer match, the buffer manager 72 will: 1) decrement the multicopy cache's frame pointer if the copy number is ">1" or 2) delete the multicopy cache's frame pointer/copy number entry and place the frame pointer into the reclaim queue 98 if the copy number is "1".

If the buffer manager 72 does not find a matching frame pointer, the buffer manager 72 searches the frame pointer's buffer header in external memory 36 (see Figure 9) for the copy number. If the copy number in memory is "1", the buffer manager 72 places the frame pointer into the reclaim queue 98. If the copy number in memory is ">1", the buffer manager 72 places the frame pointer with this copy number into the multicopy cache 96, then decrements the copy number.

The buffer manager 72 constantly services the reclaim queue 98 by reading frame pointers, then "walking" the linked-list chain to return buffers to the free buffer pool 104. This activity only returns buffers for frames that had null port vectors and were queued to the reclaim queue by the port vector FIFO 70, or frames with a multicopy forwarding vector and have completed transmissions of all of the copies. Buffers linked for unicopy frames are returned directly to the free buffer pool 104 when the frame is transmitted, as described above.

If the port vector FIFO 70 is not able to place a frame pointer for a unicopy forwarding vector into an output queue 74, because that output queue 74 and its overflow area 110 in external memory 36 are full, the frame is discarded. The frame pointer is returned to the reclaim queue 98 and the discarding of the frame is noted by the management resources of the switch. If the port vector FIFO 70 is not able to place one or more frame pointers for a multicopy forwarding vector, because one or more of the output queues 74 and their overflow areas 110 in external memory 36 are full, the frame is only forwarded to the output queues with available space and the copy number placed into the multicopy queue 90 will only reflect the successfully placed frame pointers. The non-placement of the frame pointer is noted by the switch management resources for each of the ports for which the frame pointer could not be queued. If the port vector FIFO 70 is not able to place any of the frame pointers for a multicopy forwarding vector, because all of the output queues 74 and their overflow areas 110 in external memory 36 are full, the frame pointer is passed to the reclaim queue 98, and the switch management resources are duly notified.

The multicopy queue 90 is a high priority queue used by the buffer manager 72 to keep track of how many transmissions must be completed of a particular multicopy frame before all buffers (i.e., address pointers) used to store the frame can be returned to the free buffer pool 104. The write side 92 and read side 94 of this output queue hold 64 and 16 entries, respectively. The multicopy queue 90 feeds the multicopy cache 96, which is used by the buffer manager 72 to determine when to reclaim buffers. The multicopy queue internal structure is depicted in Figure 12.

The port vector FIFO 70 places a copy of a frame's frame pointer and copy number which is ">1", based on the number of frame pointers it successfully placed in the output queues 74, into the multicopy queue 90. If a particular port's output queue 74 is full, the port vector FIFO 70 cannot place a copy of the frame pointer into the output queue 74; hence it cannot include this as a successful event in determining the copy number.

Each time the buffer manager 72 reads an output queue frame pointer and finds the unicopy bit is "0", (i.e., a multicopy), it checks the multicopy cache for the frame pointer with a copy number of "1", which indicates that this is the last transmission. If this match is found, the buffer manager 72 transmits the frame and reclaims the buffers in the same manner as in the unicopy transmission, by providing the obsolete buffers to the free buffer pool 104 after the transmission of the contents of each buffer. If the match is not found, the buffer manager 72 transmits the multicopy frame and places a copy of the frame pointer with a copy number of "-1"

into the multicopy queue 90. When a host has finished using a multicopy frame pointer for a frame which was queued to the expansion bus output queue 75 or the management port output queue 77 (through the PCI interface 26), the host writes a copy of the frame pointer with a copy number of "-1" into the multicopy queue through a frame pointer register. This register is one of the registers depicted in the block of registers 60 in Figure 2.

Similar to the output queues 74, the multicopy queue 90 is structured with an input path and an output path. The input path, or write side, allows the port vector FIFO 70 and buffer manager to place frame pointers/copy numbers into the multicopy queue 90. The output path, or read side, allows the multicopy queue 90 to place frame pointers/copy numbers into the multicopy cache 96. Additional storage for frame pointers/copy numbers, termed the multicopy queue overflow 124, is provided in external memory 36.

When frame pointers/copy numbers are written into an empty multicopy queue 90, they pass from the write side 92 to the read side 94 until the read side 94 is full. Additional frame pointers/copy numbers written to the write side 92 of the multicopy queue 90 are placed into the multicopy queue overflow area 124 in external memory 36. Once the read side 94 of the multicopy queue 90 and its overflow area 124 are full, additional frame pointers/copy numbers placed into the multicopy queue begin to fill the write side 92.

The ordering of the frame pointers passing through the multicopy queue 90 is maintained, such that when space clears in the multicopy queue read side 94, frame pointers/copy numbers are moved from the multicopy queue overflow area 124 to the multicopy queue read side 94 and from the multicopy queue write side 92 to the multicopy queue overflow area 124.

The multicopy cache 96 is similar to the multicopy queue 90 but provides a searchable region for scanning frame pointers/copy numbers. The multicopy cache 96 holds up to 256 entries. The buffer manager 72 reads a frame pointer from the multicopy queue 90 and either places it into the multicopy cache 96 or processes it, depending on whether the copy number is ">1" or "-1".

In addition, each time the buffer manager 72 reads a frame pointer from the read side 78 of an output queue 74, the buffer manager 72 schedules the transmission. If the unicopy bit is "0" (meaning a multicopy frame), the buffer manager 72 scans the multicopy cache 96 for the frame pointer with a copy number of "1", which indicates this is the last transmission of this frame. If there is a match, the buffer manager 72 removes the entry and returns buffers to the free buffer pool during frame transmission. If there is not a match, the buffer manager 72 places the frame pointer with a copy number of "-1" into the multicopy queue 90 upon completion of the transmission.

Periodically, the buffer manager 72 services the multicopy queue 90 by reading a frame pointer/copy number and placing it into the multicopy cache 96 or processing it. This is done independently of frame transmission. Two cases follow depending on whether the buffer manager reads a frame pointer with a copy number ">1" or "-1":

- 1) The buffer manager 72 reads a frame pointer with a copy number ">1" from the multicopy queue 90. Provided there is room in the multicopy cache 96, it writes a new entry. If the multicopy cache 96 is full, the buffer manager 72 must clear space in the cache 96. This is done by reading one of the older frame pointer/copy numbers from the multicopy cache 96, updating that frame pointer's buffer header in external memory 36 with the copy number in the multicopy cache 96, then deleting this cache entry. Once there is space, the new frame pointer/copy number is written into the multicopy cache 96.

2) The buffer manager 72 reads a frame pointer with a copy number "-1" from the multicopy queue 90. The buffer manager 72 searches the multicopy cache 96 for a matching frame pointer with a copy number ">=1". Two cases follow, depending on whether or not the buffer manager 72 finds a frame pointer match in the multicopy cache 96:

a) The buffer manager 72 finds a frame pointer match. If the multicopy cache 96 entry's copy number is "1", the buffer manager 72 deletes the multicopy cache entry and places the frame pointer in the reclaim queue 98. If the cache entry's copy number is ">1", the buffer manager 72 decrements the copy number by "1".

b) The buffer manager 72 does not find a frame pointer match in the multicopy cache 96. This means that the matching frame pointer has previously been moved to the buffer header of the frame's linked-list chain in external memory 36. The buffer manager 72 must go to the buffer header and read the copy number. If this value (in memory) is "1", the frame is no longer needed and the buffer manager 72 places the frame pointer in the reclaim queue 98. If this value (in memory) is ">1", the buffer manager 72 places a copy of the frame pointer/copy number (that was in external memory 36) into the multicopy cache 96 and decrements the copy number by "1". If the multicopy cache 96 is full, the buffer manager clears space by moving one of the older frame pointer/copy numbers into external memory 36.

The reclaim queue 98 holds the frame pointers to the linked-list chains which are no longer needed. The buffer manager 72 writes a frame pointer to the reclaim queue when it services the multicopy cache and discovers that a frame pointer's copy number is "1" (i.e., the last transmission of the frame has completed successfully). Additionally, the port vector FIFO 70 writes frame pointers to the reclaim queue 98 under the following conditions: 1) a frame pointer's port vector is null or 2) the frame pointer could not be queued because all of the forwarding vector's output queues were full. Finally, the host writes a frame pointer to the reclaim queue 98 (using a frame pointer register) when it has finished using a unicopy frame which was queued to the expansion bus output queue 77 or the management port output queue 75.

When the buffer manager 72 processes reclaim queue entries, it walks a frame pointer's linked-list chain to return each buffer to the free buffer pool 104. The internal structure of the reclaim queue structure is not depicted, but contains only the frame pointers (14 bits) in the exemplary embodiment of the invention. The reclaim queue write side 100 holds 64 entries and the reclaim queue read side 102 side holds 16 entries.

Similar to the output queues 74, the reclaim queue 98 is structured with an input path and an output path. The input path, or write side 100, allows the buffer manager 72 to place frame pointers in the reclaim queue 98. The output path, or read side 102, allows the buffer manager 72 to read a frame pointer and return all associated buffers to the free buffer pool 104. Additional storage for frame pointers is provided in the reclaim queue overflow area 122 provided in external memory 36.

When frame pointers are written into an empty reclaim queue 98, they pass from the write side 100 to the read side 102 until the read side 102 is full. Additional frame pointers written to the write side 100 of the reclaim queue 98 are placed into the reclaim queue overflow area 122 in external memory 36. Once the read side 102 and overflow area 122 of the reclaim queue 98 are full, additional frame pointers placed into the reclaim queue 98 begin to fill the write side 100.

Figure 11 depicts an exemplary embodiment of the internal structure of the free buffer pool 104. The free buffer pool 104 is a FIFO that contains address pointers to all free buffers 140 in external memory 36. When

frames are received, the buffer manager 72 captures available address pointers from the free buffer pool 104 to store incoming data. The buffer manager 72 also allocates address pointers from the free buffer pool 104 to the host processor 28 (when requested). The host can request or return address pointers to the free buffer pool 104 by reading or writing a free buffer pool register among the registers 60 in direct input/output space. The write side 106 and the read side 108 of the free buffer pool 104 each holds 64 entries in an exemplary embodiment of the invention.

The free buffer pool 104 is structured with an input path and an output path (similar to the output queues 74). The input path, or write side 106, allows the buffer manager 72 or host 28 to place address pointers into the free buffer pool 104. The output path, or read side 108 of the free buffer pool 104 allows the buffer manager 72 to provide address pointers to the host 28 or pull address pointers from the pool 104 for storing receive frame data. Additional storage for available address pointers, the free buffer pool overflow area 120, is provided in external memory 36, as described earlier.

Upon start-up of the switch 12, the free buffer pool generates address pointers from the read side 108. As frames come in, the free list in the free buffer pool 104 is read. If there are not enough buffer pointers in the write side 106 to handle the traffic demands, the overflow area 120 is accessed to obtain more buffer pointers.

Certain embodiments of the present invention provide an advantageous arrangement and method of providing the buffer pointers upon start-up of the switch 12. When the switch 12 first powers up, it is not required for the overflow area 120 in external memory 36 to contain buffer pointers. Instead, the buffer pointers are created on the fly. The switch 12 on power up could generate and place into the overflow area 120 the buffer pointers, but there may be 16,000 or 32,000 such pointers, and this would slow up the powering on procedure of the switch 12. The present invention takes advantage of the fact that on power up, all of the buffers are free, and the identities of these buffers are known. Therefore, the buffer pointers are generated as they are needed after power up, using a counter 180, as depicted in Figure 10.

The free list count generator 180 is connected to the input of a multiplexer 182. Since the free list in the free buffer pool 104 is empty on startup, the free list counter 180 generates the buffer pointers. Once the free list reaches the highest count, it will not generate any more buffer pointers.

When a frame packet is received in the switch 12, the frame packet is broken up into fixed length buffers. Typically, frames vary in size. The buffers are 256 bytes in size and the data portion of a buffer is 240 bytes. Following transmission of the contents of a buffer, the buffer pointers are put into the reclaim queue 98 or, if the buffer chain can be walked, directly into the free list of the free buffer pool 104. During operation of the switch 12, any address pointers returned to the free buffer pool 104 pass from the write side 106 to the read side 108. If the read side 108 becomes full, additional address pointers are passed to the overflow area 120. Once the read side 108 and the overflow area 120 are full, additional address pointers placed into the free buffer pool 104 will begin to fill the write side 106 of the pool 104 again.

Figure 13 depicts a schematic representation of the internal arrangement of the multicopy cache 96 in accordance with an embodiment of the present invention. As briefly discussed earlier, the time order of the entries to the multicopy cache 96 is maintained. In the present invention, this maintaining of a time order is not done by time stamping, as in the prior art, but by physical ordering in a memory. The multicopy cache 96 of the present invention also avoids the use of validity bits, and instead encodes validity, as will be discussed.

Referring to Figure 13, the multicopy cache 96 is configured as a four-way set-associative memory. An entry into the multicopy cache 96 includes a frame pointer and its copy number, as explained earlier. The lowest six bits of the frame pointer determine the row in the set-associative cache 96 in which the entry will be stored. In the illustrated embodiment of the invention, there are sixty-four rows in the cache 96, although other numbers of rows are not limited if the cache size is made larger.

The set-associative cache 96 is divided into four columns, each of which can be searched in parallel. When the buffer manager 72 stores an entry into the cache 96, the entry always enters the first column, the uppermost (51:39) bits of the row indicated by the six least significant bits of the frame pointer. The row is read, all of the entries are shifted to the right by 13 bits, and the row is written back. The entry that is actually written into the cache 96 includes the upper eight bits of the frame pointer that form an address tag, and the five-bit copy number associated with the frame pointer. When the entry is read out of the cache 96, the frame pointer is re-formed with the address tag and the bits that index the row number of the cache 96.

The oldest entry in the cache 96 is removed from the cache 96 if the row is full and a new entry to the row is written. As described earlier with respect to the buffer headers 142, the copy number associated with the frame pointer that is removed is written into the buffer header 142 of the frame in external memory pointed to by the removed frame pointer. Hence, the frames (i.e., the buffers 140) stored in external memory 36 serve as an overflow area for the multicopy cache 96 to store copy numbers.

One of the advantageous features of the present invention is that there is no separate valid bit in the set-associative cache 96. When the copy number is 00000, then the buffer manager 72 knows that the entry is no longer valid and removes the entry from the cache 96. This simplifies the organization of the cache. Another advantage of the cache 96 of the present invention is that it allows a very fast search to be performed, since the buffer manager 72 needs only to examine a single row, already determined by the frame pointer that has exited the multicopy queue 90. The four entries in the row are examined in parallel, further increasing the speed of the search. Although described as a four-way set-associative memory, this is exemplary only as the memory can be n-way set-associative without departing from the scope of the invention.

From the above description, it should be understood that the present invention maintains a time order (age) of the cache entries by physical positioning of the entries in the cache on a row basis. In other words, the physical position of the entry in the cache is an indication of the relative age of an entry. The aging of an entry is performed by the physical re-ordering of the entries in the memory.

Certain embodiments of the present invention provide customization of the latency of frames switched by the switch 12, on a port by port basis. Referring to Figure 14, the port vector FIFO 70 examines the programmed switch mode of the receive port to determine when to place the frame pointer and associated information into the appropriate output queue 74 of the transmit port. For a first mode (low latency mode), the port vector FIFO 70 has no restrictions on when to place the frame pointer into the output queue(s) 74. For a second mode (intermediate latency mode), the port vector FIFO 70 places the frame pointer into the output queue(s) 74 only after 64 bytes of the frame have been received. For a third mode (high latency mode), the port vector FIFO 70 places the frame pointer into the output queue(s) 74 only after the frame has been completely received.

There are some special cases which alter the timing of when the port vector FIFO 70 passes frame pointers to the output queues 74: 1) frames forwarding from a first or second mode 10 Mb/s port to a 100 Mb/s port; 2) frames forwarding to the management port 30 and 3) frames forwarding to the expansion bus port. In case 1), the 10 Mb/s port to 100 Mb/s port rate mismatch forces the forwarding mode to be the third, high latency mode. In case 2), all frames passed to the management port are third mode frames. In case 3), any frame forwarding to the expansion bus port uses the switch mode of the expansion bus port 26. When a multicopy port vector contains one of the special case ports, the queuing of the frame pointers for the entire port vector becomes that of the longest latency switch mode represented in the port vector. For example, if a frame is received by a first or a second mode port, but its multicopy forwarding port vector contains the management port 30, the switch mode is the third mode. In this situation, a copy of the frame pointer is placed into all the output queues 74 only after the frame has been completely received.

The switch modes will now be described in more detail. The switch mode that applies to the input (i.e., receive) port determines forwarding latency (how soon the switch 12 will forward a frame once it begins receiving the frame) and the ability to reduce fragment/error propagation to output ports. The second, intermediate latency mode is the default for each port; however, the switch mode is programmable on a per-port basis in the registers 60.

In all of the three modes, frame data received at the receive FIFO 52 of an internal MAC port is forwarded to a buffer 140 in the external memory 52 as soon as possible. At approximately the same time, the rules checker 42 or 58, receives the destination address and source address, the receive port number, the frame pointer and some additional information, then performs the appropriate lookup. Once the lookup is completed, the rules checker 42 or 58 returns the frame pointer and the forwarding port vector to the port vector FIFO 70.

The port vector FIFO 70 places the frame pointer in the write side 76 of the output queues 74 for the output port(s) identified in the port vector. The receive port's switch mode defines the latency between when the port vector FIFO 70 receives the port vector (and the frame pointer) and places the frame pointer into the output queue(s) 74. This is described for the three modes below. Once the frame pointer passes to the read side 78 of the output queues 74, the buffer manager 72 reads the frame pointer and schedules transmission. The buffer manager 72 begins moving frame data from the address specified by the frame pointer. Once the transmit FIFO 54 of the MAC port has been primed to its start point (and assuming the medium is available for transmission of data), frame transmission commences.

The first mode is designed to provide the lowest latency. Frames are received and forwarded at line-rate speed. In this first mode, there is no network error protection because a frame is queued for transmission before it can be determined whether the frame is a fragment (i.e., <64 bytes in length) or contains a CRC error. In the first mode, frame reception may not complete before frame transmission at the output port(s) commences. If a receive frame terminates in a runt or with an invalid CRC, the receive MAC marks the buffer header 142 in external memory 36 to indicate these conditions. The transmit MAC guarantees that if transmission commences on a frame which later terminates as a runt or with an invalid CRC, the MAC will generate a bad CRC. If the transmit MAC has not started a frame's transmission and the buffer header 142 indicates the frame terminated in a runt or an invalid CRC, the buffer manager 72 will not forward the frame to the output port.

The second mode provides low latency for forwarding frames and some network error protection. Frames are received and forwarded after sixty-four or more bytes have been received. This allows the switch 12 to filter (i.e., not forward) fragments of frame; however, it does not completely filter CRC error frames that are greater than sixty-four bytes.

In the second mode, frame pointers for frames which achieve the sixty-four byte threshold at the receive MAC are queued to the appropriate output queue(s) 74. Frames which fail to achieve the minimum sixty-four byte threshold are deleted and their frame pointers are not placed in output queue(s) 74. If a receive frame greater than or equal to sixty-four bytes terminates with an invalid CRC, the receive MAC marks the buffer header 142 in external memory 36 to indicate this condition. If transmission has commenced on a frame greater than or equal to sixty-four bytes which later terminates with an invalid CRC, the transmit MAC will complete the transmission with a bad CRC. If the transmit MAC has not started a frame transmission and the buffer header 142 indicates the frame (greater than or equal to sixty-four bytes) terminated in an invalid CRC, the buffer manager 72 returns the frame pointer to the reclaim queue 98 (for a unicopy forward) or the multicopy queue 96 (for a multicopy forward) without forwarding to the output port(s) 74.

The third mode is a store-and-forward mode that provides the highest level of network error protection among the three modes, but has a higher forwarding latency. Frames are completely received before the switch 12 forwarding them to output ports. In this mode, the switch 12 screens out all fragments and CRC error frames before forwarding. In the third mode, once a valid frame completes successfully at the receiver (i.e., greater than or equal to sixty-four bytes with a valid CRC), the frame pointer is queued to the appropriate output queue(s) 74. Frames which terminate in a receive error (invalid CRC, runt (>64 bytes) etc.) are deleted and their frame pointers are not placed in output queue(s) 74.

The port vector FIFO 70 makes the decision to put the port vector into an output queue 74, in dependence on the selected mode of the receive port and the amount of data that has been received. In the embodiment described above, there are three thresholds, although there are different numbers of thresholds in other embodiments. In the exemplary embodiment, these thresholds are: 1) receiving n bytes (e.g. 6 bytes) where $n < 64$ bytes; 2) receiving 64 bytes; and 3) receiving all of the frame.

The present invention forwards frames to the output queues 74 based on thresholds. The port vector FIFO 70 re-orders the sequence of transmission based on amount of type of data received and the mode in which the port is programmed. Although the exemplary embodiment makes forwarding decisions based on the amount of received data, other embodiments of the invention make forwarding decisions based on other factors, such as the types of data received.

In implementing the forwarding scheme of the present invention, the buffer manager 72 maintains a table 160 in a cache memory (CAM) 161 that associates a frame pointer with a receive port. Every time the port vector FIFO 70 receives a new port vector and frame pointer from the rules checker 42 or 58, it makes an association to determine whether the receive port has finished receiving a frame, and if not, how much of the frame has already been received. The port vector FIFO 70 does not receive any information regarding the identity of the receive port from the rules checker 42 or 58. The only information the port vector receives that provides any identification of the port are the frame pointers.

The port vector FIFO 70 queries the address table 160 with the frame pointer. Either the address table returns the receive port if the frame is still being received, or the address table 160 cannot find the frame pointer which means that the frame has already been received. Once the frame is completely received, the frame pointer is moved out of the address table 160. This means that the third threshold (the frame complete) is met. Accordingly, the frame pointer may be dropped into the output queue 74 immediately.

If the address table 160 returns the receive port, the port vector FIFO 70 puts the frame pointer and associated information into a holding area 162 and begins monitoring two signals from that receive port. These two signals flag one of three events. The first event is flagged when the port receives n bytes. At that point, if that port is in the first mode, the port vector FIFO 70 starts processing the frame pointer by sending it to the appropriate output queue 74. If the receive port is not in the first mode, the port vector FIFO 70 waits until it receives the a signal indicating occurrence of the second event. If this port is in the second mode, then the port vector FIFO 70 releases the frame pointer from the holding area 162 to enter the proper output queue 74. Finally, if the receive port is in the third mode, then the port vector FIFO 70 awaits receipt of the flag indicating that the frame is complete. Every receive port (reference numeral 164 in Figure 14) maintains this flag, and provides this information to the port vector FIFO 70. It is up to the port vector FIFO 70 to determine the port associated with the frame pointer. The port vector FIFO 70 maintains the information identifying the mode each port is in. In summary, upon receiving a frame pointer, the port vector FIFO 70 first queries the address table 160 of the buffer manager 72 to determine the receive port, determines the mode for that receive port, and then monitors the flags from that receive port and releases the frame pointer according to the mode and the flags.

The queue structure of the present invention has been described as incorporated into a multiport switch, as an example. However, the queue architecture and method of the present invention is applicable to other types of technologies, as will be apparent to those of ordinary skill in the art. For example, the queue architecture may be used in microprocessors or other digital processors that have on-chip queues. The invention is therefore not limited to multiport switches, such as those described above.

Although the present invention has been described and illustrated in detail, it is to be clearly understood that the same is by way of illustration and example only and is not to be taken by way of limitation, the spirit and scope of the present invention being limited only by the terms of the appended claims.

CLAIMS

WHAT IS CLAIMED IS:

1. A queue that queues and prioritizes entries, comprising:
 - a queue write side forming a queue input at which entries to the queue are received, each entry having a priority level;
 - a plurality of queue read sides forming a queue output at which entries from the queue write side are received, queued and output from the queue, the plurality of queue read sides having different levels of queuing priority from a highest priority queue read side to a lowest priority queue read side; and
 - control logic configured to examine entries to the queue write side to determine the priority level of the entries and to direct the entries from the queue write side to one of the queue read sides based on the priority levels of the entries.
2. The queue of Claim 1, wherein the number of priority levels of the entries is equal to the number of queue read sides with different priority levels.
3. The queue of Claim 2, wherein each entry has priority bits in which the priority of that entry is encoded, the control logic being configured to examine the priority bits of each entry to the queue.
4. The queue of Claim 3, wherein the queue write side has at least twice the queuing capacity of any of the individual queue read sides.
5. The queue of Claim 4, wherein the queuing capacity of the highest priority queue read side is less than the queuing capacity of the lowest priority queue read side.
6. The queue of Claim 5, wherein the output queue read sides are configured to permit detection of presence of entries in the output queue read sides and control of emptying of the individual output queue read sides such that the output queue read sides are emptied in the order of highest to lowest output queue read side priority, so that a higher priority output queue read side has no queued entries before entries in a lower priority queue read side are output.
7. The queue of Claim 6, further comprising an overflow area for storing the entries, the overflow area coupled to receive the entries from the queue write side and return the entries to one of the queue read sides.
8. The queue of Claim 7, wherein the control logic is configured to direct entries having a priority level below a first specific priority level from the queue write side to the overflow area when the queue read sides with priority levels matching the priority levels of the entries have no available queuing capacity, and to direct entries from the overflow area to the matching priority level queue read sides when the matching priority level queue read sides regain available queuing capacity.

9. The queue of Claim 8, wherein the control logic is further configured to discard entries from the output queue write side having a priority level above a second specific priority level when the queue read sides with priority levels matching the priority levels of the entries have no available queuing capacity.

10. The queue of Claim 9, wherein the entries are frame pointers that point to locations in memory where frames of data are located.

11. A multiport network switch for a packet switched network, comprising:

a plurality of ports through which packets from the network are received and from which packets to the network are transmitted;

a plurality of queues corresponding to the plurality of ports, at least one of the queues having a queue write side forming a queue input at which entries to the queue are received, each entry having a priority level, and a plurality of queue read sides forming a queue output at which entries from the queue write side are received, queued and output from the queue, the plurality of queue read sides having different levels of queuing priority from a highest priority queue read side to a lowest priority queue read side; and

control logic configured to examine entries to the queue write side to determine the priority level of the entries and to direct the entries from the queue write side to one of the queue read sides based on the priority levels of the entries.

12. The switch of Claim 11, wherein the number of priority levels of the entries is equal to the number of queue read sides with different priority levels.

13. The switch of Claim 12, wherein each entry has priority bits in which the priority of that entry is encoded, the control logic being configured to examine the priority bits of each entry to the queue.

14. The switch of Claim 13, wherein the queue write side has at least twice the queuing capacity of any of the individual queue read sides.

15. The switch of Claim 14, wherein the queuing capacity of the highest priority queue read side is less than the queuing capacity of the lowest priority queue read side.

16. The switch of Claim 15, wherein the output queue read sides are configured to permit detection of presence of entries in the output queue read sides and control of emptying of the individual output queue read sides such that the output queue read sides are emptied in the order of highest to lowest output queue read side priority, so that a higher priority output queue read side has no queued entries before entries in a lower priority queue read side are output.

17. The switch of Claim 16, further comprising an interface to an overflow area for storing the entries, the interface providing the overflow area with entries from the queue write side and return the entries from the overflow area to one of the queue read sides.

18. The switch of Claim 17, wherein the control logic is configured to direct entries having a priority level below a first specific priority level from the queue write side through the interface to the overflow area when the queue read sides with priority levels matching the priority levels of the entries have no available queuing capacity, and to direct entries through the interface from the overflow area to the matching priority level queue read sides when the matching priority level queue read sides regain available queuing capacity.

19. The switch of Claim 18, wherein the control logic is further configured to discard entries from the output queue write side having a priority level above a second specific priority level when the queue read sides with priority levels matching the priority levels of the entries have no available queuing capacity.

20. The switch of Claim 19, wherein the entries are frame pointers that point to locations in memory where frames of data are located.

21. A method of queuing and prioritizing entries comprising the steps of:
receiving entries to the queue in a queue write side that forms an input of the queue;
determining a level of priority of the entries to the queue write side;
directing entries exiting the queue write side to one of a plurality of queue read sides of the queue, based on the determined level of priority of the entries, the queue read sides having different levels of queuing priority, from a highest priority queue read side to a lowest priority queue read side.

22. The method of Claim 21, further comprising determining the available capacity of the queue read sides and directing an entry from the queue write side to an overflow area when the queue read side having the same level of priority as the priority level of the entry has no available capacity, and directing the entry from the overflow area to that queue read side when that queue read side regains available capacity.

23. The method of Claim 22, further comprising discarding an entry exiting from the queue write side when the queue read side having the same level of priority as the priority level of the entry has no available capacity, if the level of priority is above a threshold level of priority.

24. The method of Claim 23, further comprising the steps of polling the queue read sides to detect presence of entries in the queue read sides and emptying the individual queue read sides in the order of highest to lowest queue read side priority, so that a higher priority queue read side has no queued entries before entries in a lower priority queue read side are output.

1/10

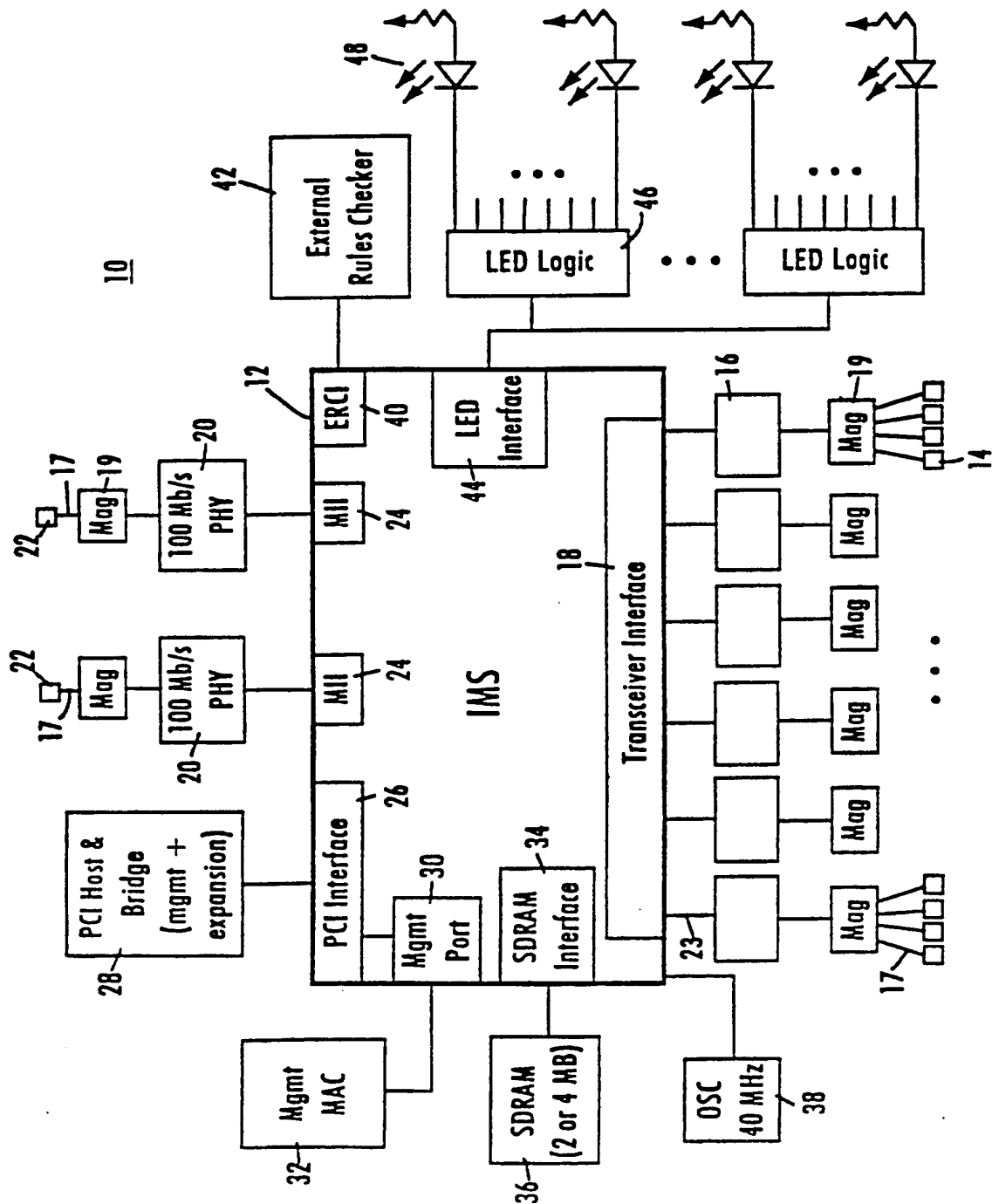
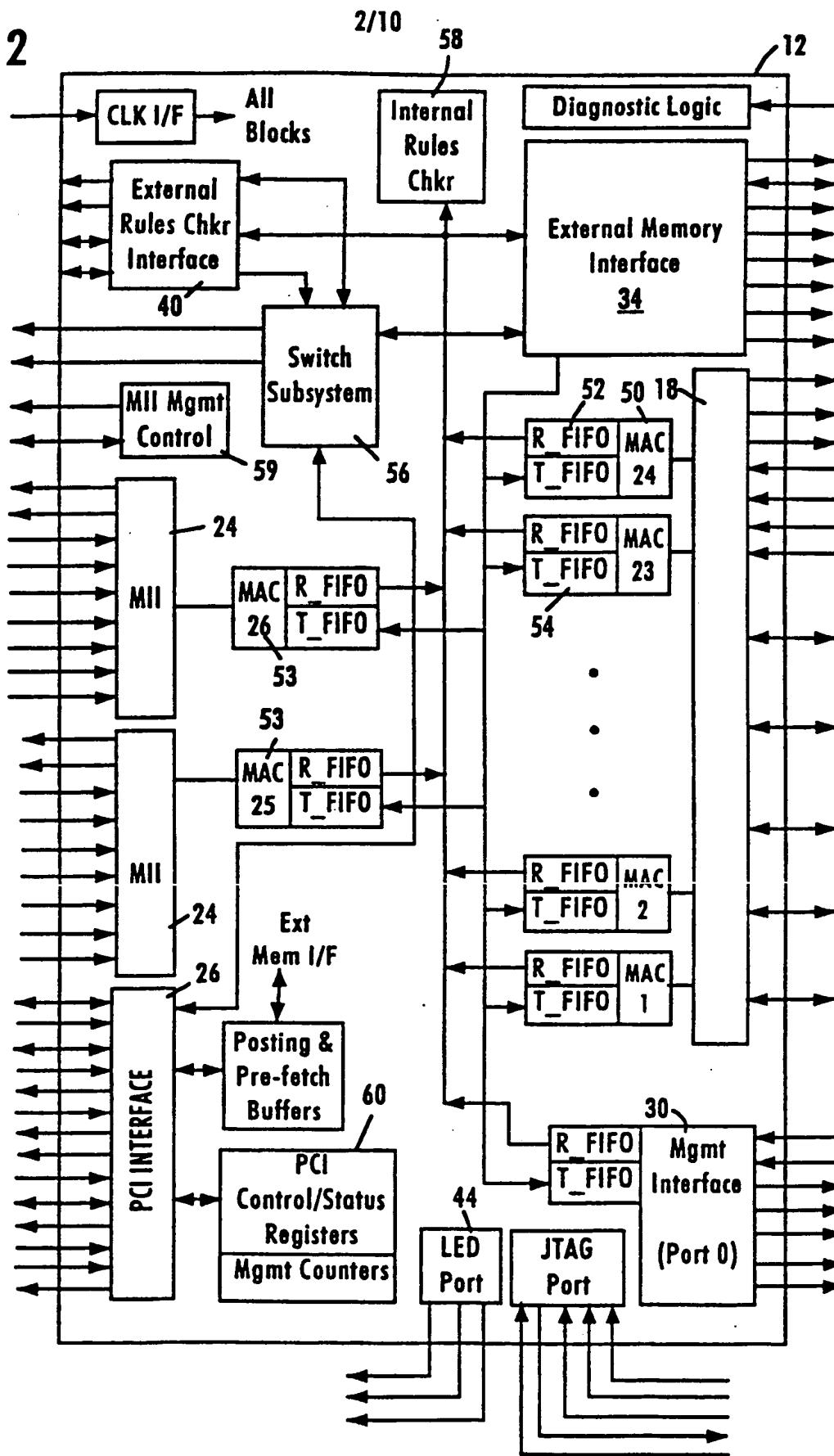


Fig. 1

Fig. 2



4/10

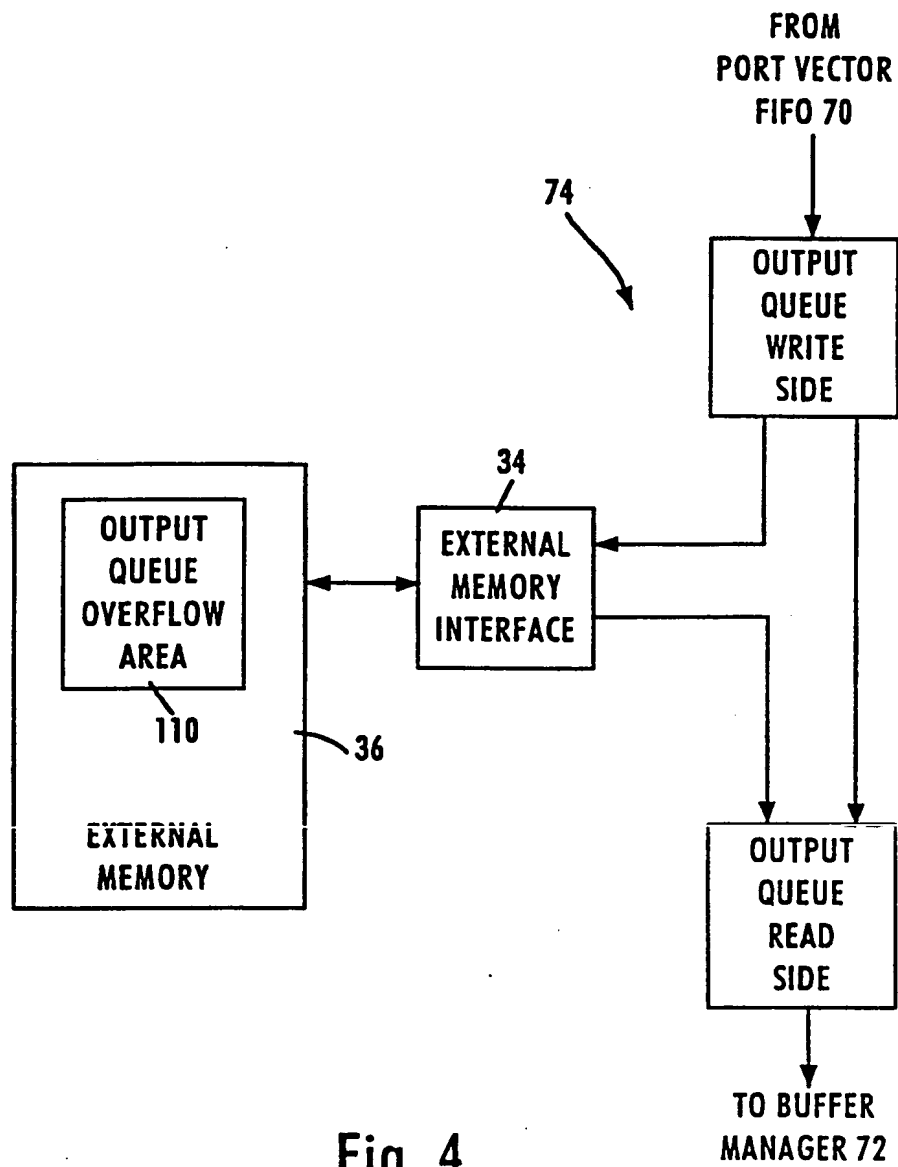


Fig. 4

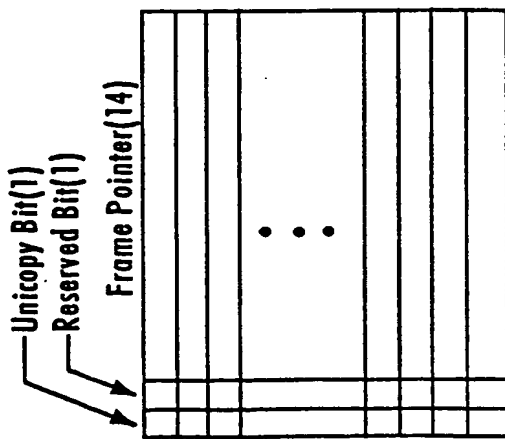


Fig. 5

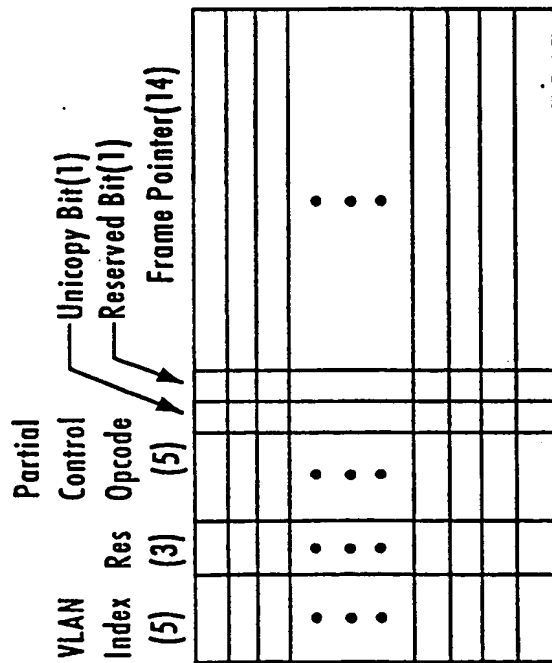


Fig. 6

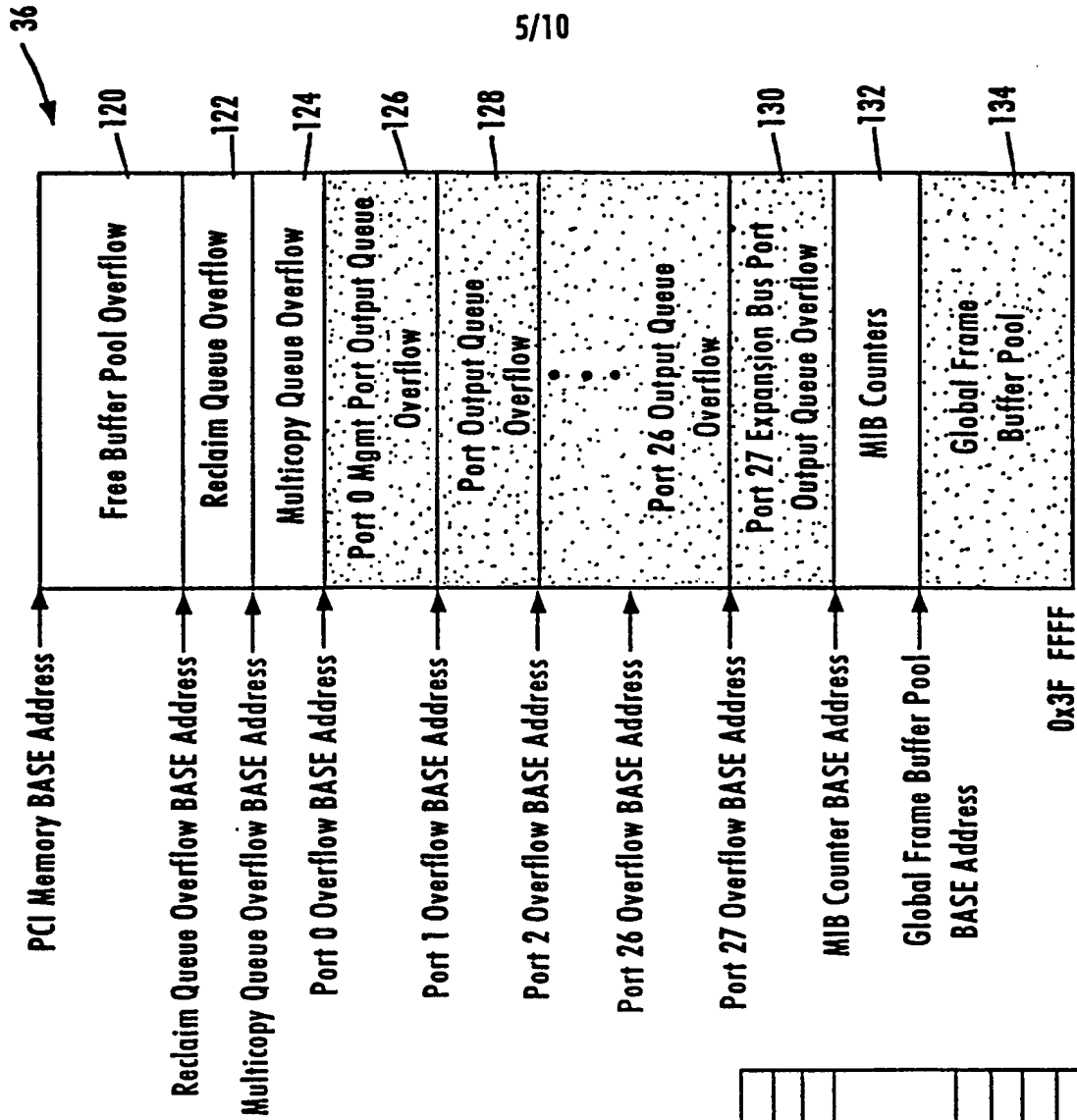


Fig. 7

6/10

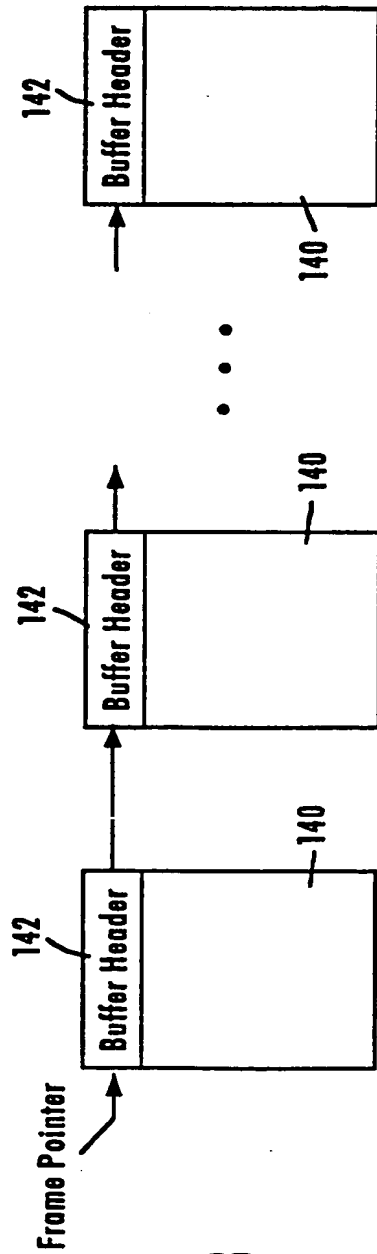


Fig. 8

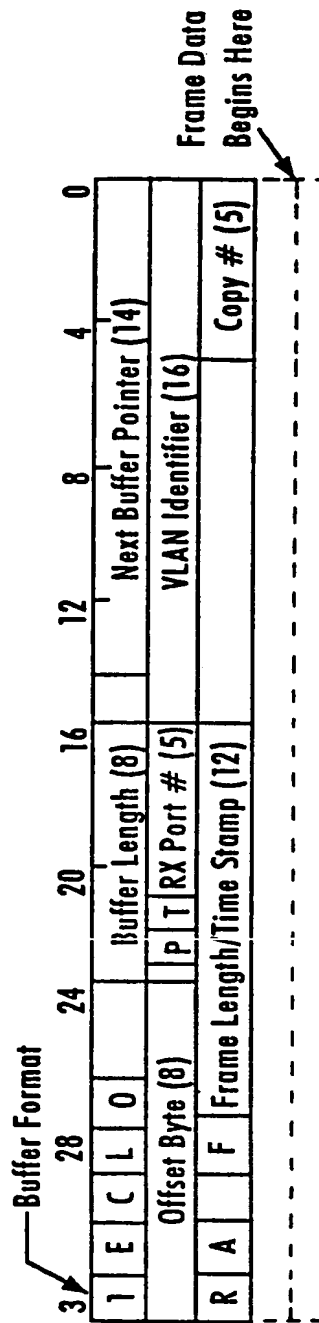


Fig. 9a

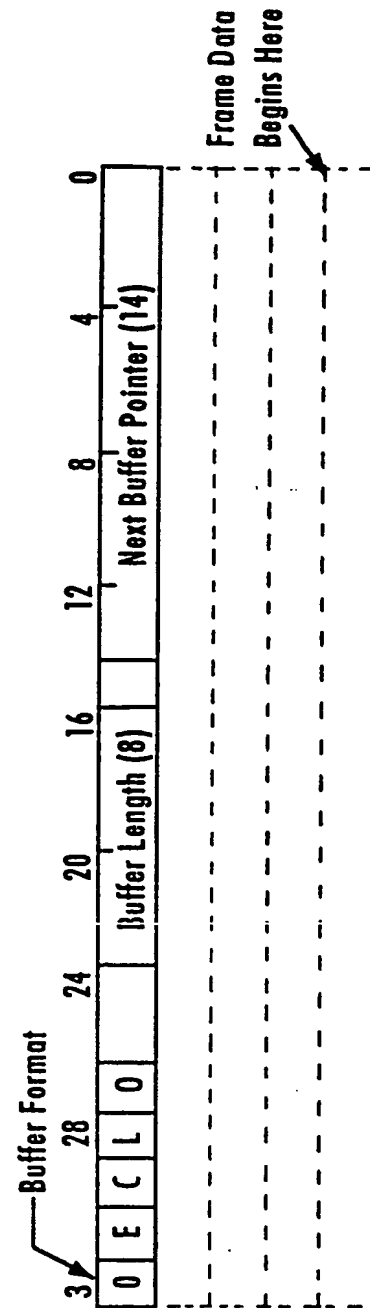


Fig. 9b

7/10

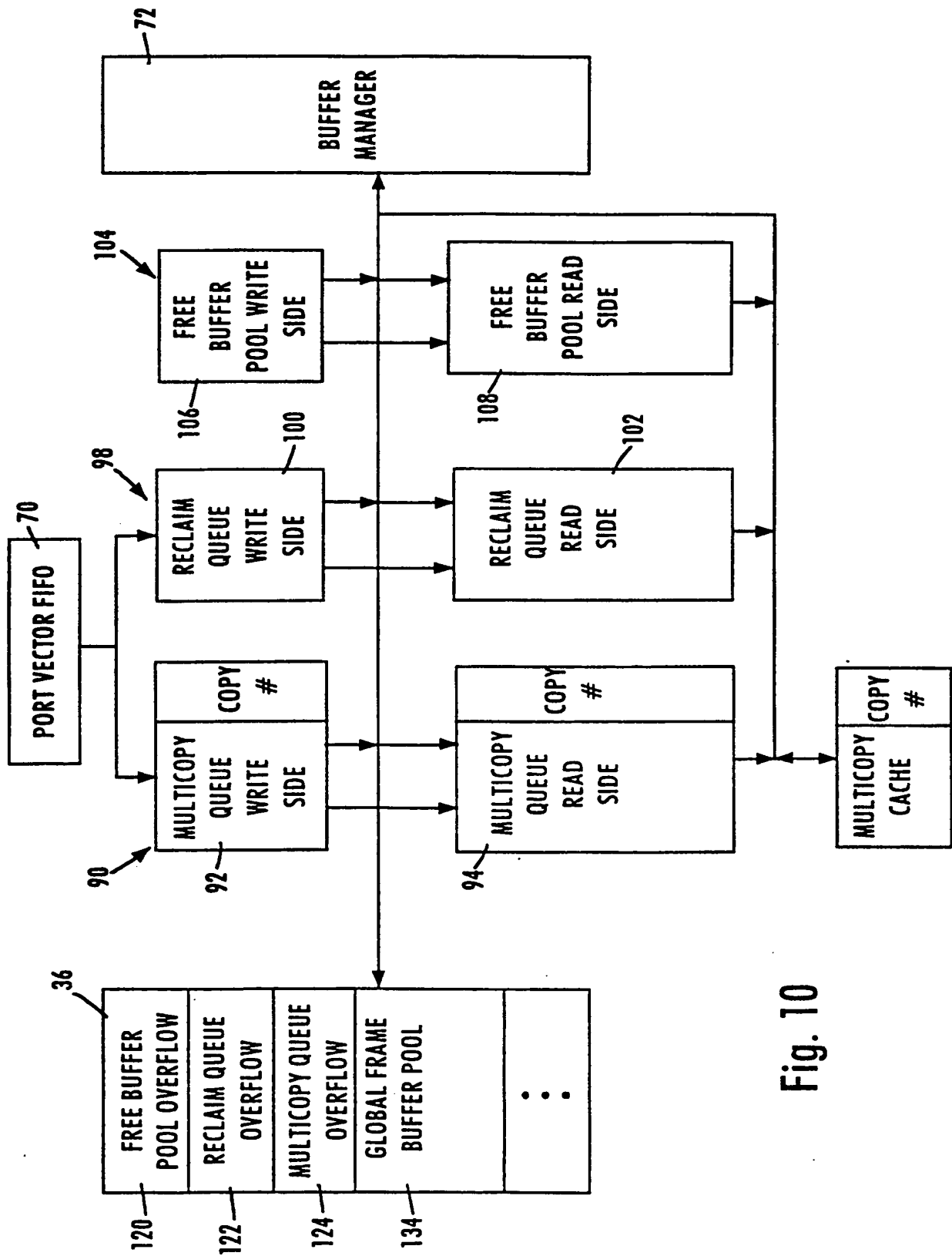


Fig. 10

8/10

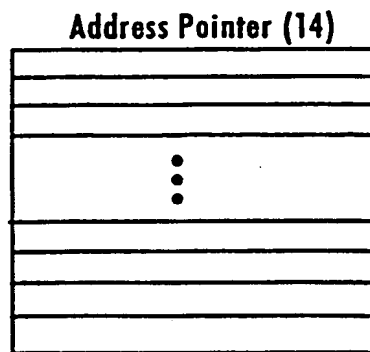


Fig. 11

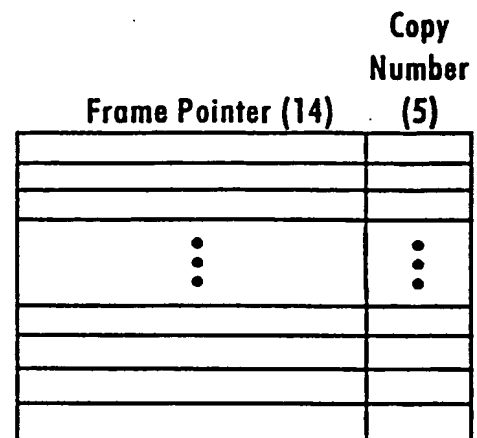


Fig. 12

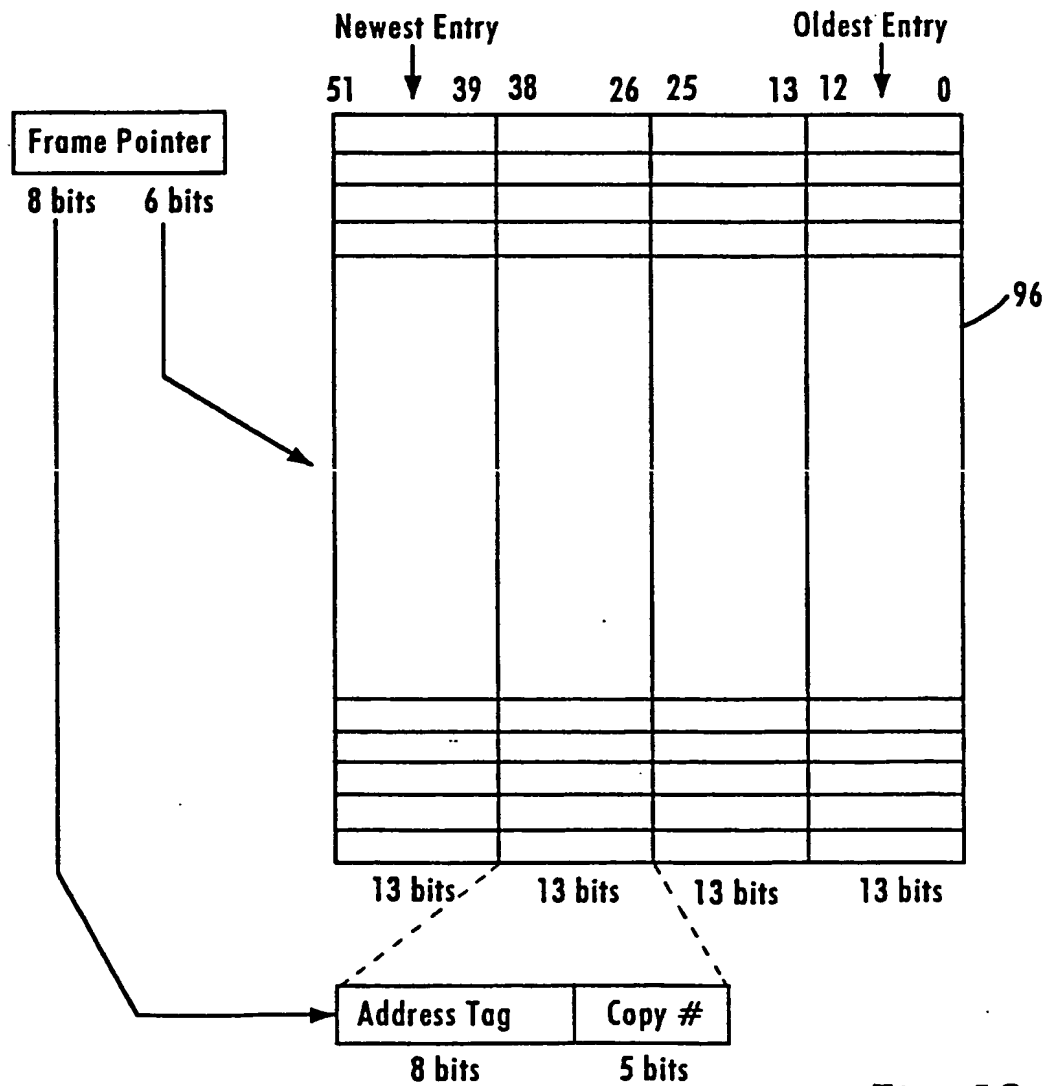


Fig. 13

9/10

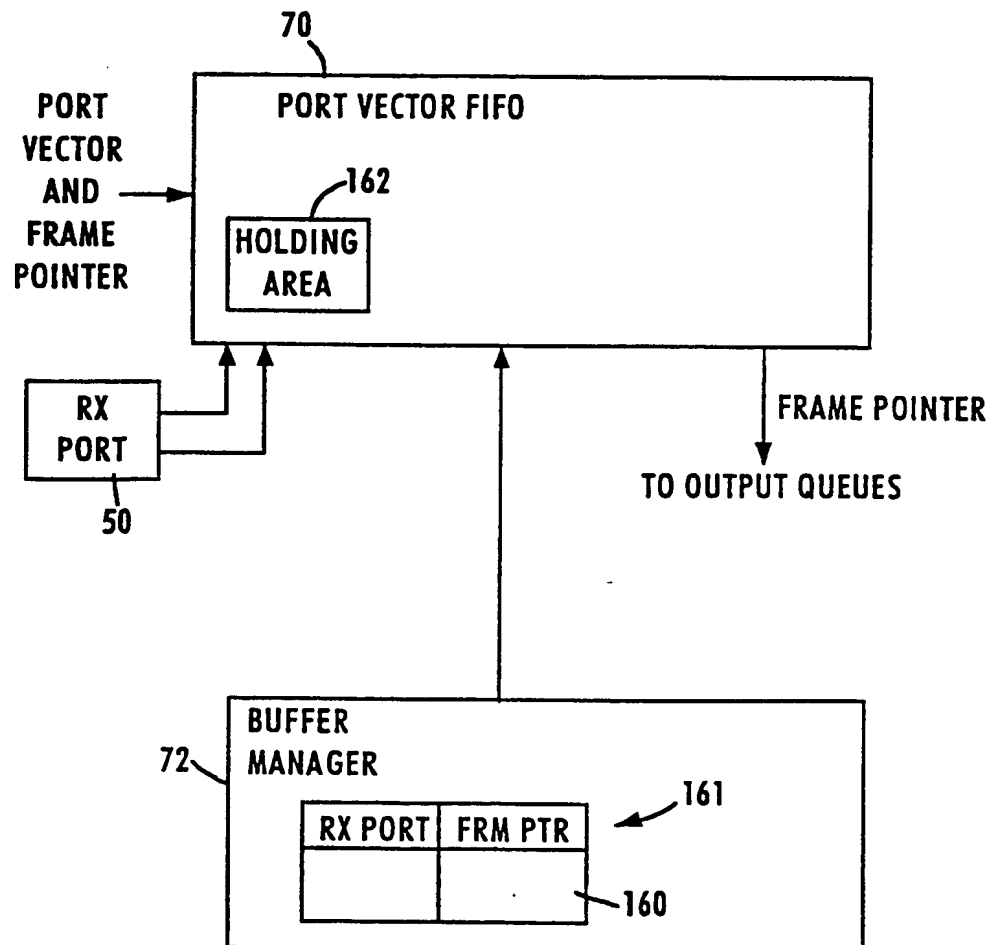


Fig. 14

10/10

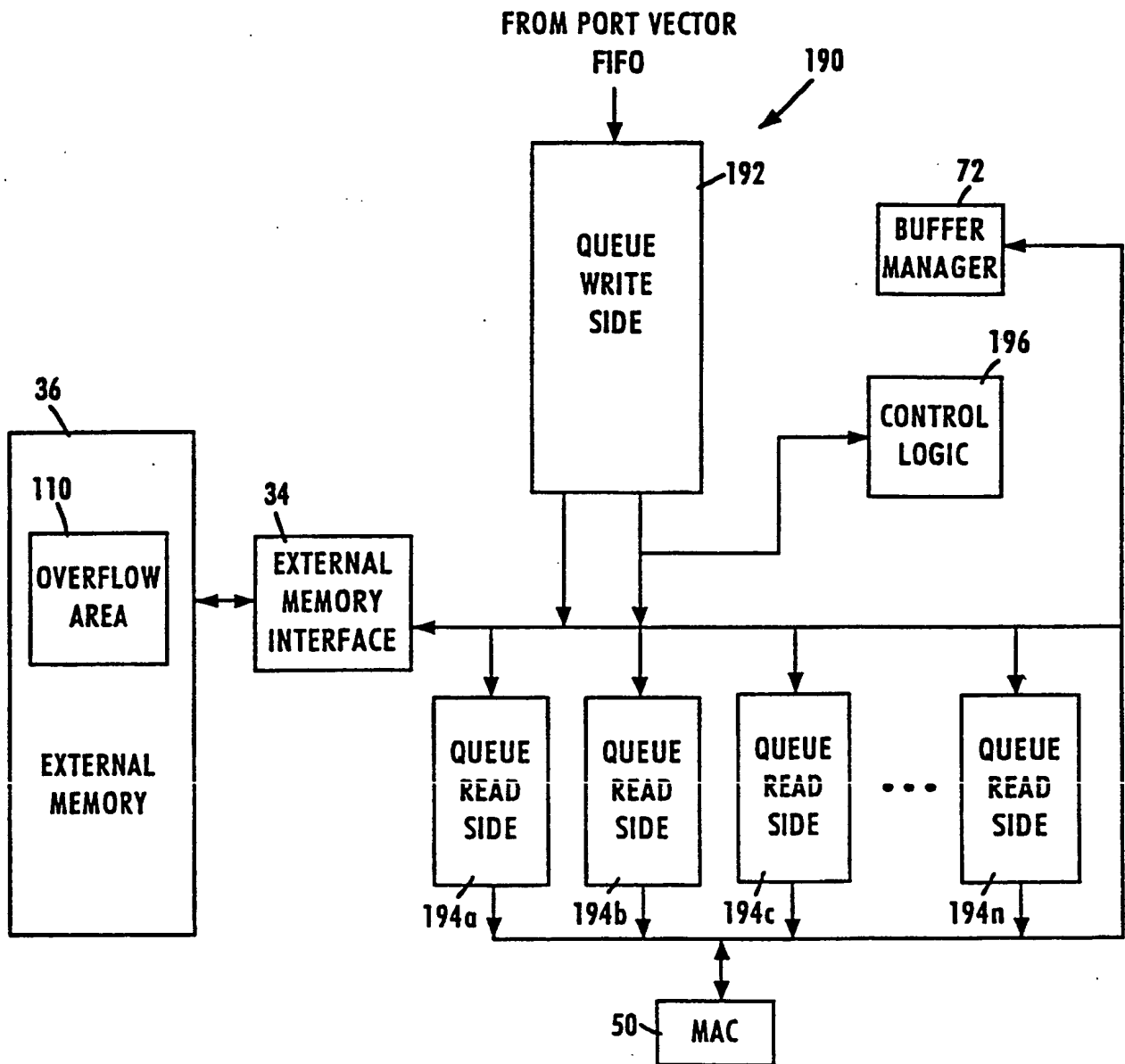
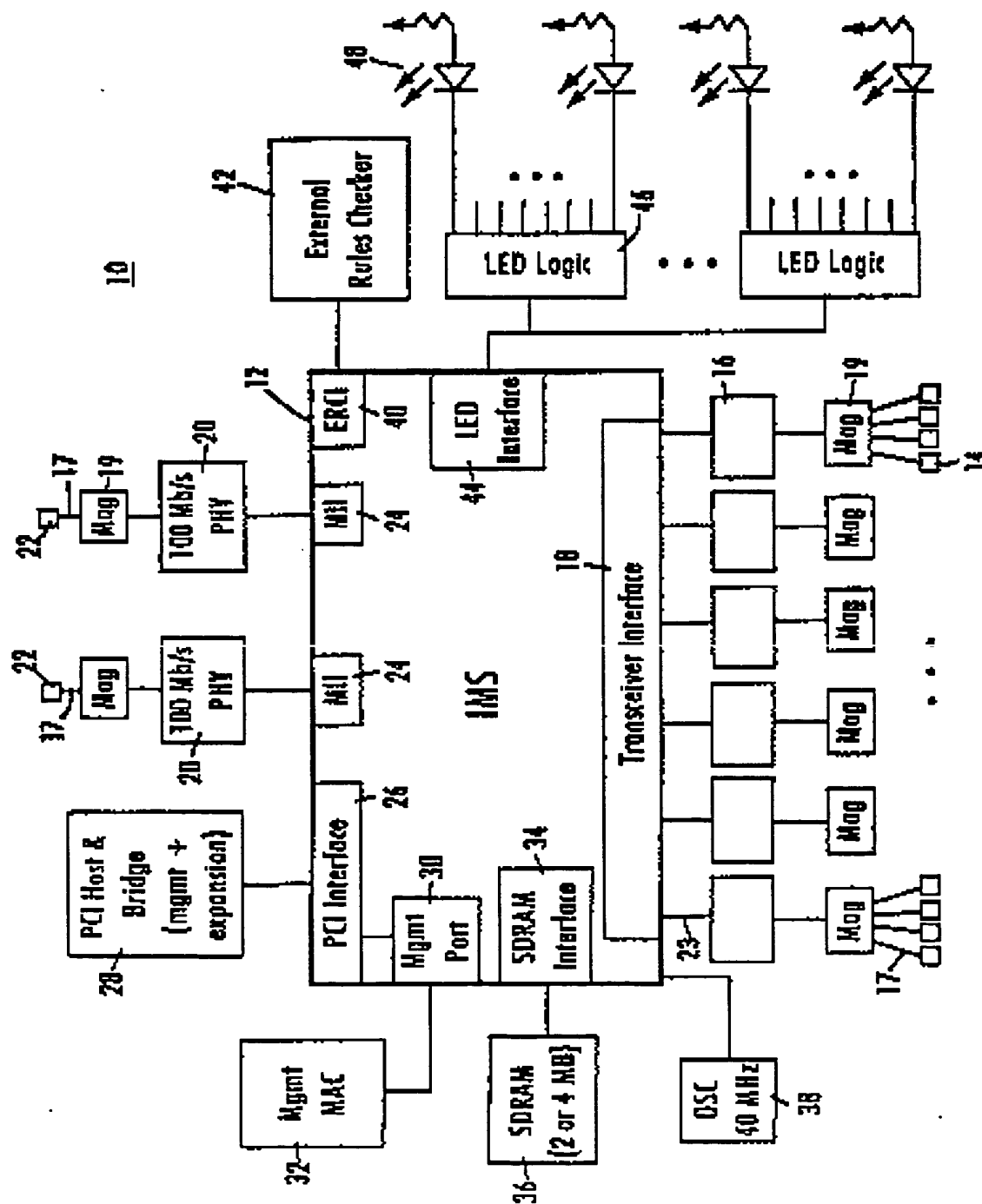


Fig. 15

THIS PAGE BLANK (USPTO)

1110



101

Fig. 2

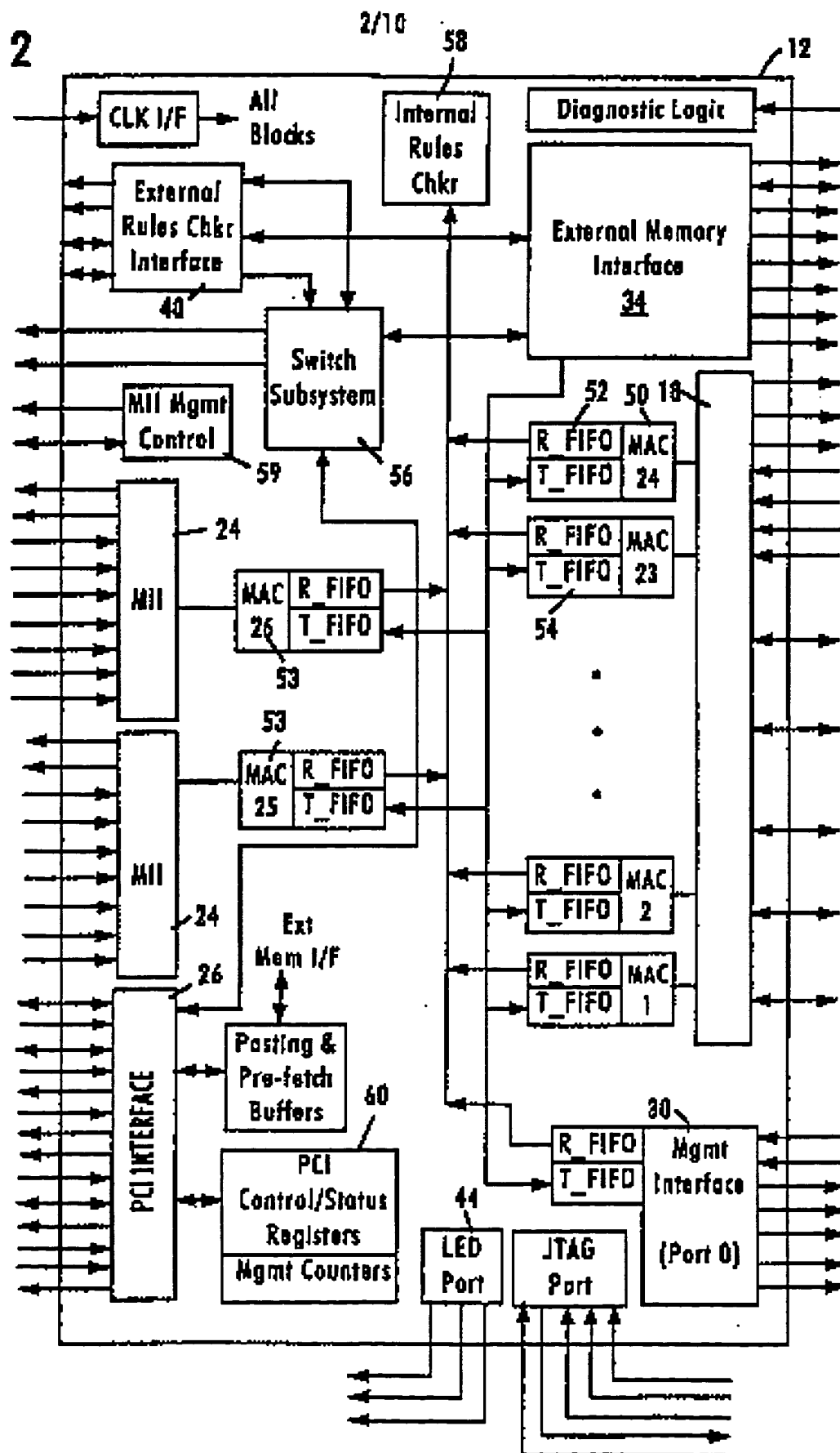
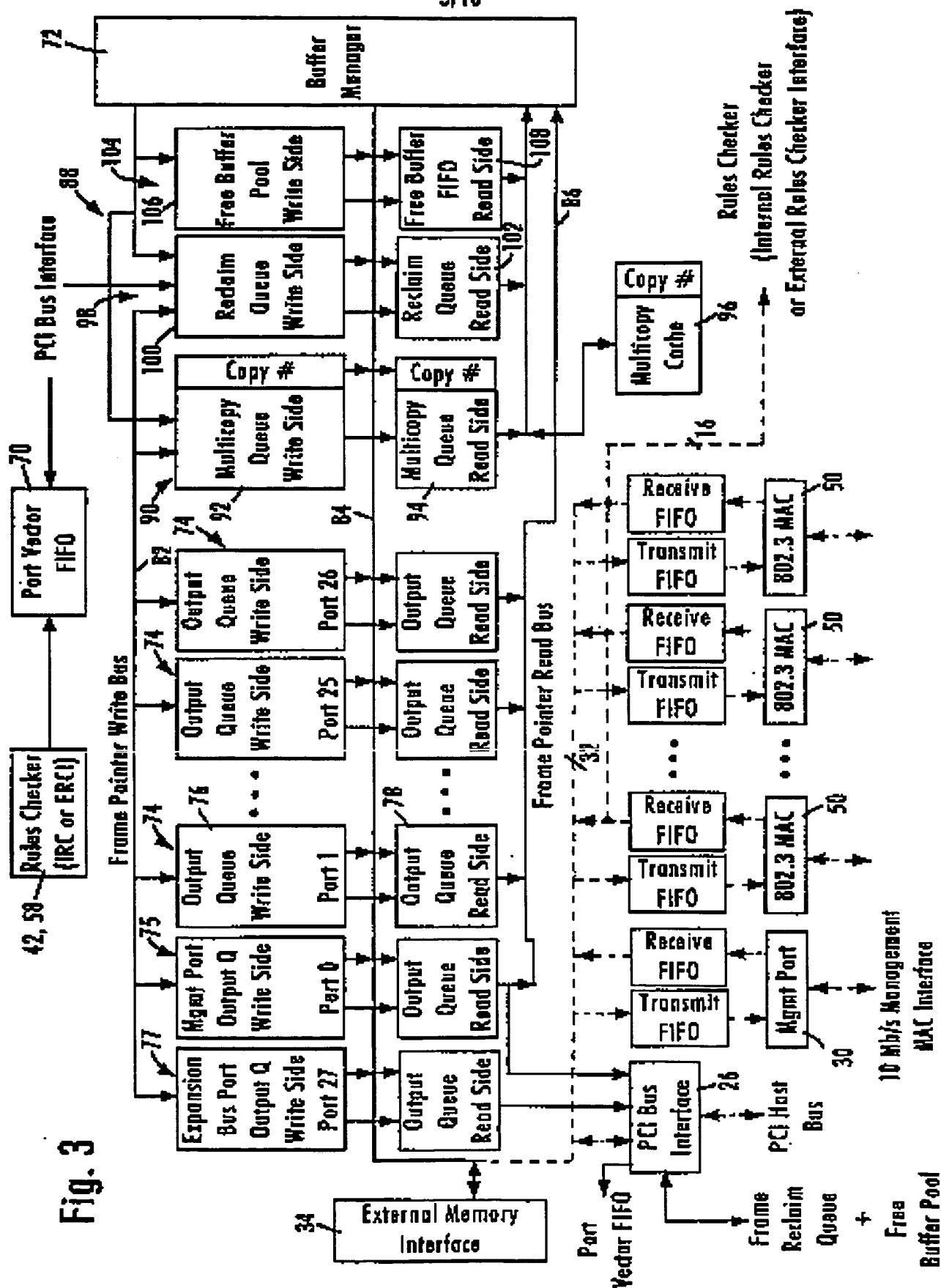


Fig. 3



4/10

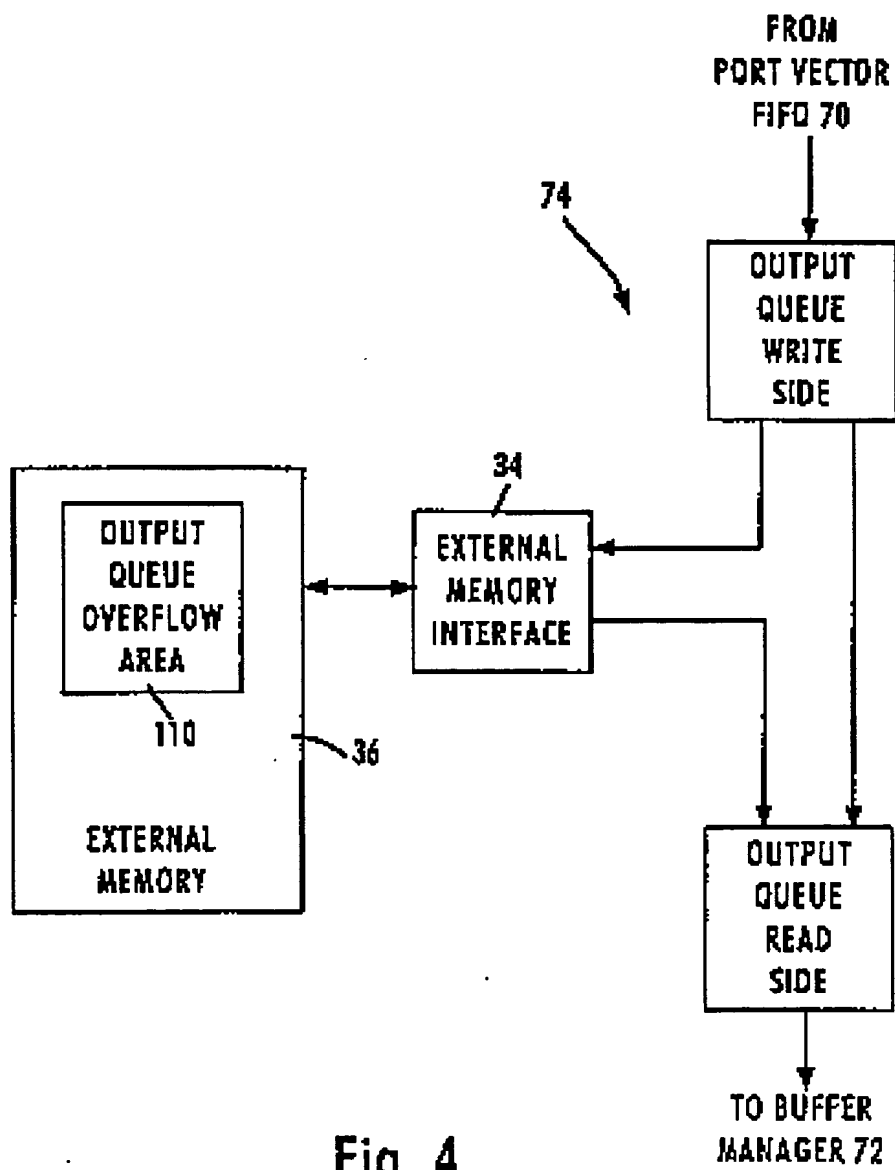


Fig. 4

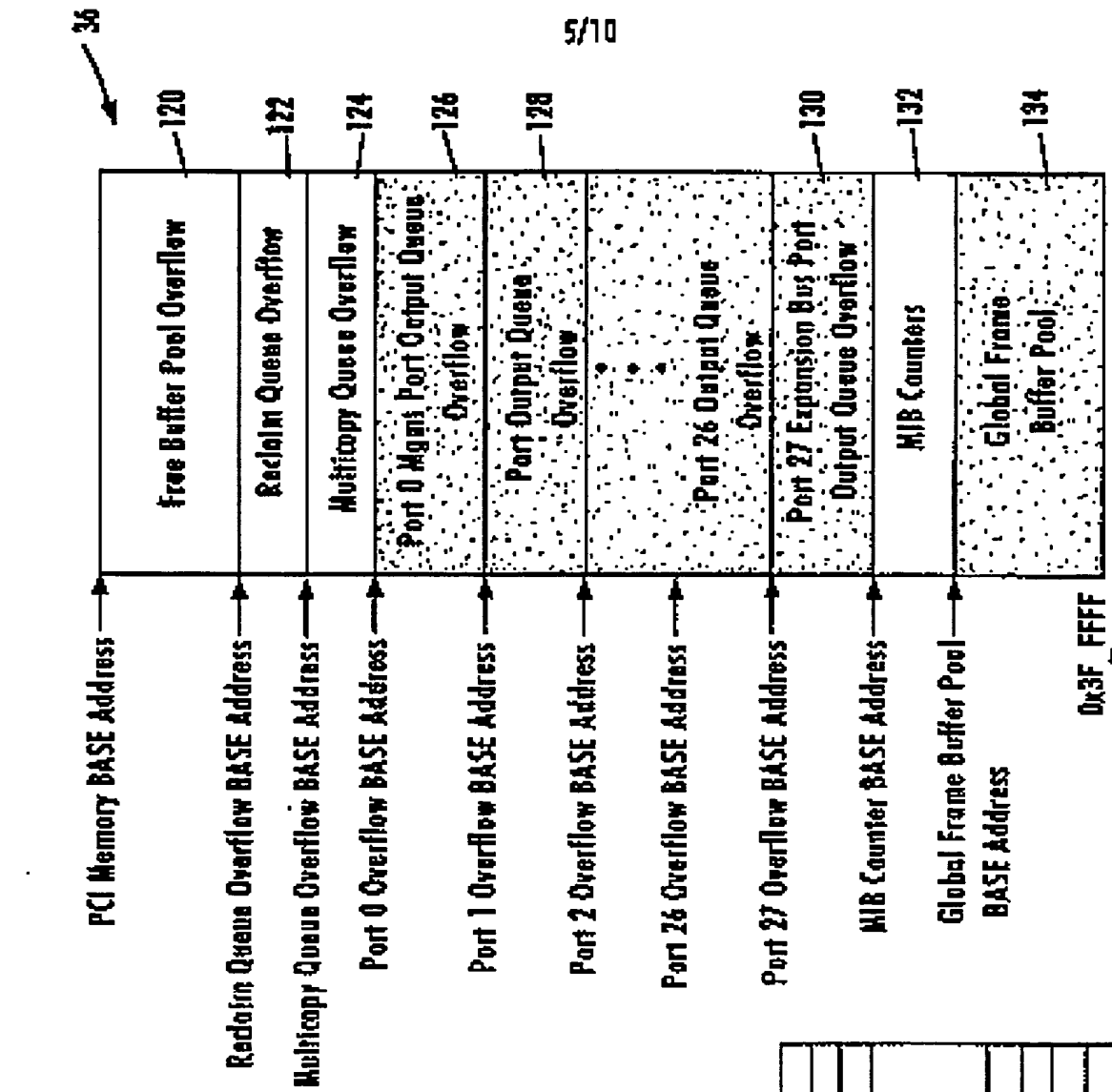


Fig. 5

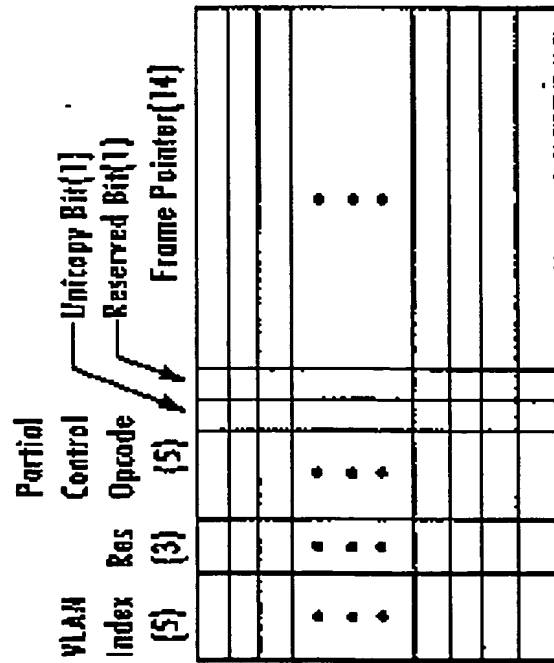


Fig. 6

Fig. 7

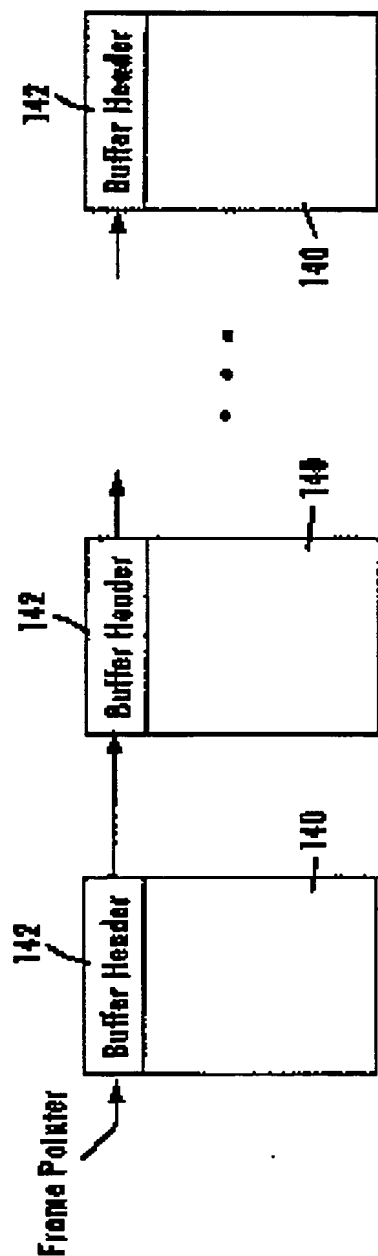


Fig. 8

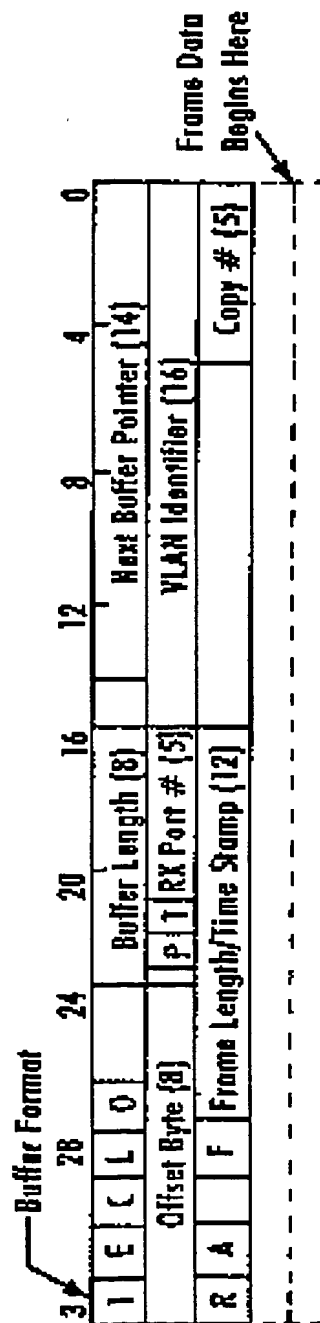


Fig. 9a

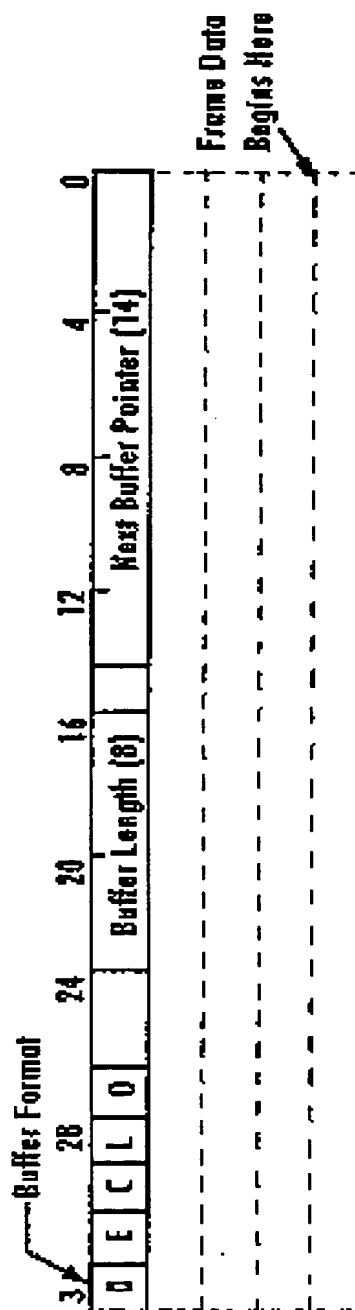


Fig. 9b

7/10

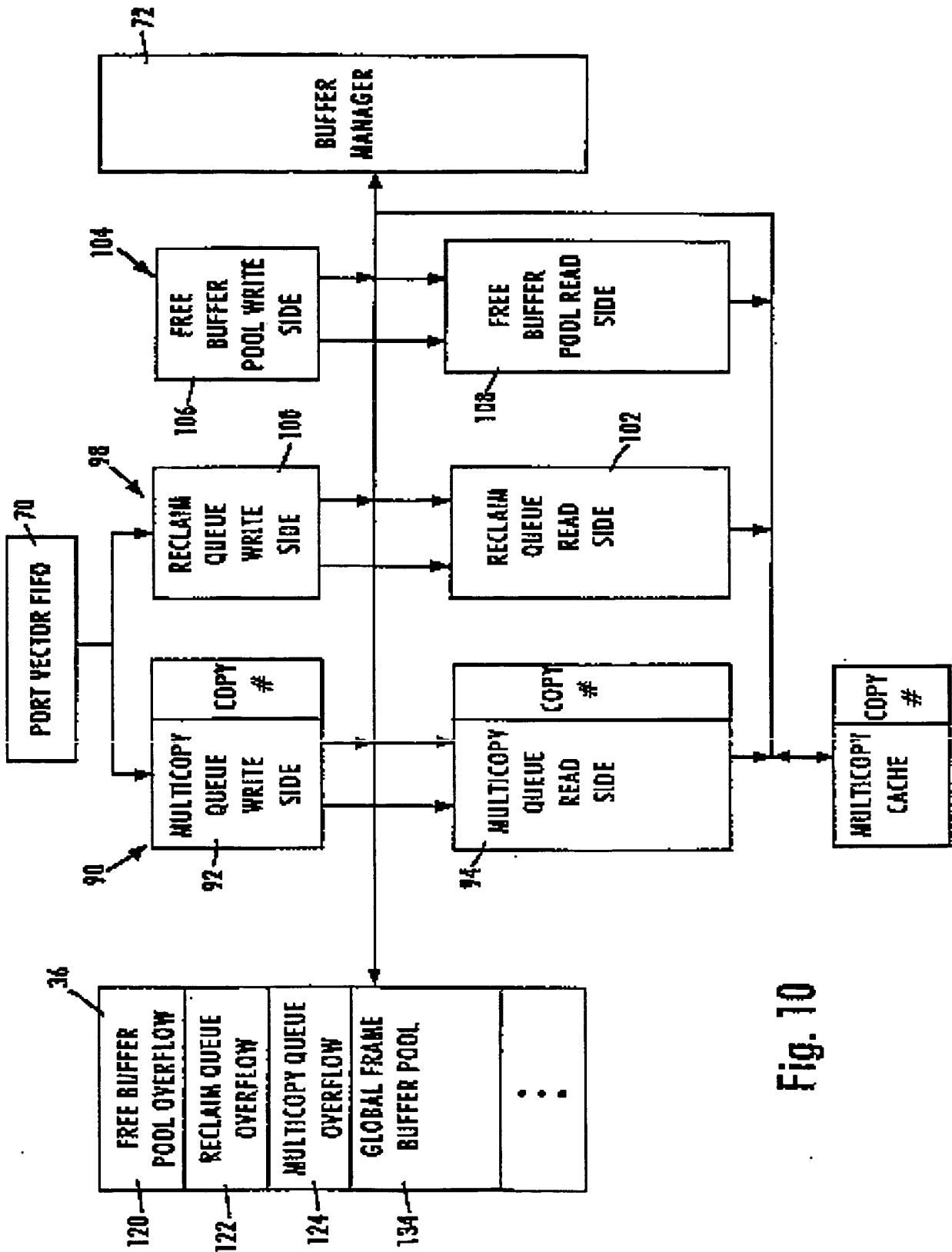


Fig. 10

9/10

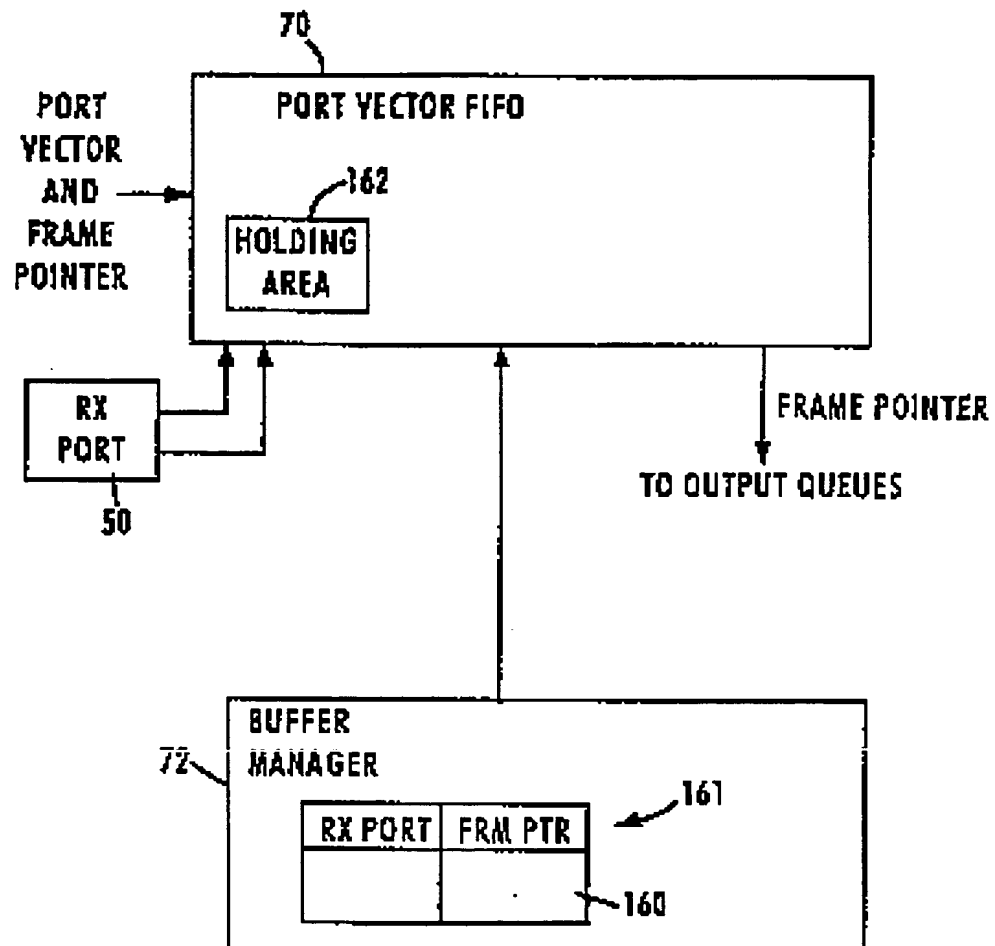
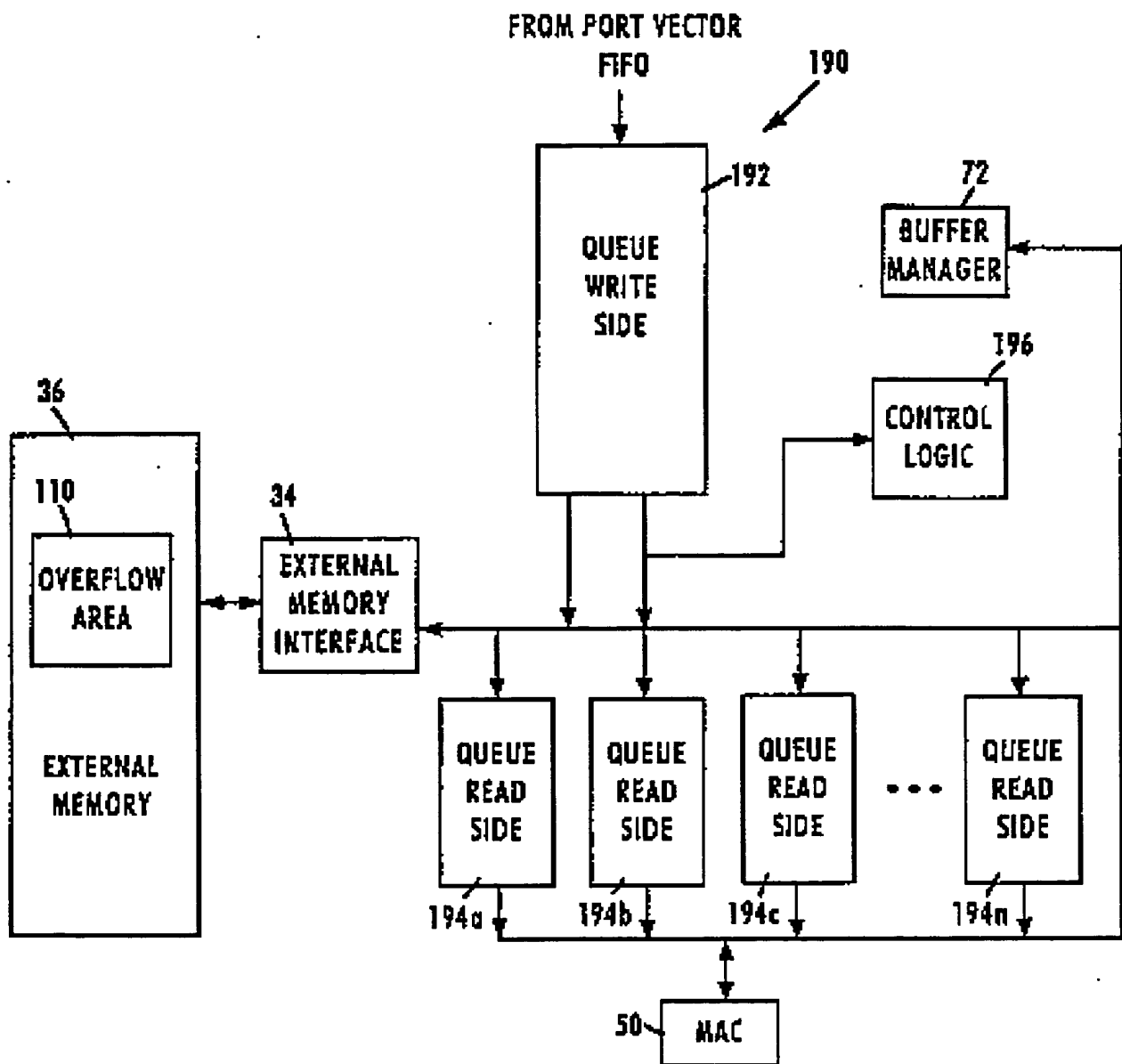


Fig. 14

10/10

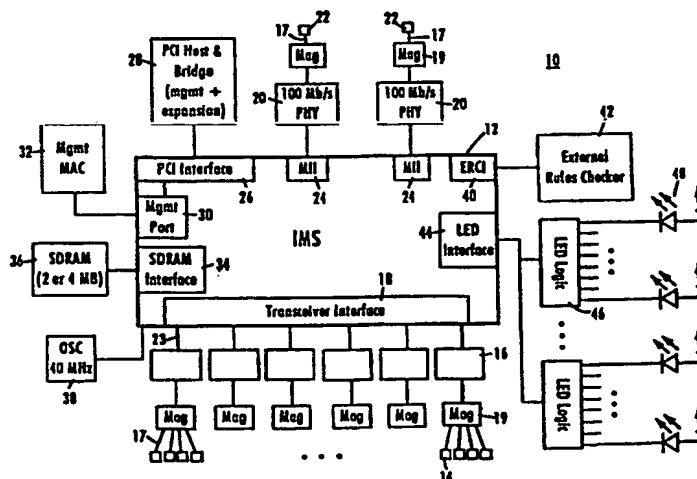




INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification ⁶ : H04L 12/44, 12/46, 29/06	A3	(11) International Publication Number: WO 98/36587 (43) International Publication Date: 20 August 1998 (20.08.98)
(21) International Application Number: PCT/US98/02357 (22) International Filing Date: 6 February 1998 (06.02.98) (30) Priority Data: 60/038,025 14 February 1997 (14.02.97) US 08/993,147 18 December 1997 (18.12.97) US (71) Applicant: ADVANCED MICRO DEVICES, INC. [US/US]; One AMD Place, Mail Stop 68, Sunnyvale, CA 94088-3453 (US). (72) Inventors: ERIMLI, Bahadir; 2101 California Street #109, Mountain View, CA 94040 (US). CRAYFORD, Ian; 5380 Eileen Drive, San Jose, CA 95129 (US). EGBERT, Chandan; 3632 Bloomsbury Way, San Jose, CA 95132 (US). (74) Agent: ZAHRT, William, D., II.; Advanced Micro Devices, Inc., One AMD Place, Mail Stop 68, Sunnyvale, CA 94088-3453 (US).	(81) Designated States: JP, KR, European patent (AT, BE, CH, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE). Published <i>With international search report.</i> <i>Before the expiration of the time limit for amending the claims and to be republished in the event of the receipt of amendments.</i> (88) Date of publication of the international search report: 14 January 1999 (14.01.99)	

(54) Title: QUEUING STRUCTURE AND METHOD FOR PRIORITIZATION OF FRAMES IN A NETWORK SWITCH



(57) Abstract

A method and arrangement for queuing data in a prioritized manner have a queue with a single queue write side in which data entries are input to the queue. The queue also has a plurality of queue read sides, with each read side having a different priority level. An entry exiting from the queue write side is examined to determine the priority level of the entry, and then placed into the queue read side with the matching priority level as the entry. The queue read sides form the output of the queue, and are polled and emptied so that the higher priority queues are emptied completely before emptying lower priority queues.

FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece			TR	Turkey
BG	Bulgaria	HU	Hungary	ML	Mali	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MN	Mongolia	UA	Ukraine
BR	Brazil	IL	Israel	MR	Mauritania	UG	Uganda
BY	Belarus	IS	Iceland	MW	Malawi	US	United States of America
CA	Canada	IT	Italy	MX	Mexico	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NE	Niger	VN	Viet Nam
CG	Congo	KE	Kenya	NL	Netherlands	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NO	Norway	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	NZ	New Zealand		
CM	Cameroon	KR	Republic of Korea	PL	Poland		
CN	China	KZ	Kazakstan	PT	Portugal		
CU	Cuba	LC	Saint Lucia	RO	Romania		
CZ	Czech Republic	LI	Liechtenstein	RU	Russian Federation		
DE	Germany	LK	Sri Lanka	SD	Sudan		
DK	Denmark	LR	Liberia	SE	Sweden		
EE	Estonia			SG	Singapore		

INTERNATIONAL SEARCH REPORT

International Application No
PCT/US 98/02357

A. CLASSIFICATION OF SUBJECT MATTER

IPC 6 H04L12/44 H04L12/46 H04L29/06

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

IPC 6 H04L H04J

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	US 4 630 254 A (TSENG CHONG-WEI) 16 December 1986 see column 1, line 24 - column 6, line 21 see column 7, line 8 - line 21 see column 9, line 55 - column 11, line 57 see column 13, line 30 - column 14, line 27	1-3, 11-13, 21
Y		4, 14, 22
A	---	5, 15
	-/--	

☒ Further documents are listed in the continuation of box C.

☒ Patent family members are listed in annex.

* Special categories of cited documents :

- "A" document defining the general state of the art which is not considered to be of particular relevance
- "E" earlier document but published on or after the international filing date
- "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)
- "O" document referring to an oral disclosure, use, exhibition or other means
- "P" document published prior to the international filing date but later than the priority date claimed

- "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
- "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
- "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.
- "&" document member of the same patent family

Date of the actual completion of the international search

11 November 1998

Date of mailing of the international search report

20. 11. 98

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,
Fax: (+31-70) 340-3016

Authorized officer

Vaskimo, K

INTERNATIONAL SEARCH REPORT

Internat' Application No
PCT/US 98/02357

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	EP 0 544 188 A (MITSUBISHI ELECTRIC CORP) 2 June 1993 see column 1, line 15 - line 54 see column 3, line 3 - column 4, line 25 see column 6, line 19 - column 7, line 22 see column 8, line 7 - line 31 see column 14, line 40 - column 15, line 15 see claims 1,6,7,10,11	1,2,11, 12,21
Y A	---	22 3,13
X	US 5 541 922 A (PYHAELAMMI SEPPÖ) 30 July 1996 see abstract see column 1, line 7 - column 2, line 60 see column 3, line 29 - line 67 see column 4, line 29 - column 5, line 2	1,2,11, 12,21
Y A	---	22 3,13
Y	YOUN CHAN JUNG ET AL: "PERFORMANCE ANALYSIS OF PACKET SWITCHES WITH INPUT AND OUTPUT BUFFERS" COMPUTER NETWORKS AND ISDN SYSTEMS, vol. 26, no. 12, 1 September 1994, pages 1559-1580, XP000458464 see abstract see page 1560, line 11 - page 1562, line 5 see page 1572, line 10 - page 1573, line 2 see figure 12	4,14
A	---	1-3,5, 11-13,15
Y	CHEN J S C ET AL: "OPTIMAL BUFFER ALLOCATION FOR PACKET SWITCHES WITH INPUT AND OUTPUT QUEUEING" COMMUNICATIONS: CONNECTING THE FUTURE, SAN DIEGO, DEC. 2 - 5, 1990, vol. VOL. 3, no. -, 2 December 1990, pages 1936-1941, XP000218904 INSTITUTE OF ELECTRICAL AND ELECTRONICS ENGINEERS see abstract see page 1937, paragraph III.A - page 1938, paragraph III.B	4,14
A	---	1-3,5, 11-13,15
	--- -/--	

INTERNATIONAL SEARCH REPORT

Internat Application No

PCT/US 98/02357

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

Category	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
L,X	US 5 774 453 A (FUKANO ET AL.) 30 June 1998 see abstract see column 1, line 6 - column 5, line 22	1-3, 11-13, 21,22
L,Y L,A	-& JP 08 288953 A1 November 1996 see the whole document & AU 50726 96 A31 October 1996 see the whole document ---	23,24 4,14
Y	EP 0 491 489 A (AMERICAN TELEPHONE & TELEGRAPH) 24 June 1992 see abstract see column 1, line 3 - column 2, line 28 see column 3, line 28 - column 4, line 44 see column 6, line 33 - column 8, line 36 see column 13, line 56 - column 16, line 37 see figure 7	22
A	---	1,11,21, 23
Y	EP 0 526 104 A (AMERICAN TELEPHONE & TELEGRAPH) 3 February 1993 see the whole document	22
A	---	1,11,21, 23
Y	US 5 079 762 A (TANABE) 7 January 1992 see column 1, line 39 - column 2, line 14 see column 4, line 5 - line 36	22
A	---	1,11,21, 23
Y	US 5 583 858 A (HANAOKA HIDETOSHI) 10 December 1996 see column 1, line 51 - column 2, line 37 see column 4, line 43 - column 5, line 40 see column 8, line 62 - column 9, line 62 see claims 1,7	23,24
A	---	1,11,21, 22
Y	US 5 313 579 A (CHAO HUNG-HSIANG J) 17 May 1994 see abstract see column 1, line 16 - column 2, line 49 see column 4, line 31 - column 6, line 19	23,24
A	---	1,11,21, 22

	-/--	

INTERNATIONAL SEARCH REPORT

Intern. J. Application No
PCT/US 98/02357

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	US 5 555 265 A (KAKUMA ET AL.) 10 September 1996 see column 3, line 18 - line 49 see column 4, line 49 - column 5, line 47 see column 8, line 50 - column 9, line 28 see column 10, line 34 - column 11, line 14 see column 11, line 62 - column 12, line 13 see column 12, line 66 - column 13, line 61 see figures 1,2,5,7,11	23,24
A	-----	1,11,21, 22

INTERNATIONAL SEARCH REPORT

International application No.
PCT/US 98/02357

Box I Observations where certain claims were found unsearchable (Continuation of item 1 of first sheet)

This International Search Report has not been established in respect of certain claims under Article 17(2)(a) for the following reasons:

1. ☐ Claims Nos.:
because they relate to subject matter not required to be searched by this Authority, namely:

2. ☐ Claims Nos.:
because they relate to parts of the International Application that do not comply with the prescribed requirements to such an extent that no meaningful International Search can be carried out, specifically:

3. ☐ Claims Nos.:
because they are dependent claims and are not drafted in accordance with the second and third sentences of Rule 6.4(a).

Box II Observations where unity of invention is lacking (Continuation of item 2 of first sheet)

This International Searching Authority found multiple inventions in this international application, as follows:

see additional sheet

1. ☒ As all required additional search fees were timely paid by the applicant, this International Search Report covers all searchable claims.

2. ☐ As all searchable claims could be searched without effort justifying an additional fee, this Authority did not invite payment of any additional fee.

3. ☐ As only some of the required additional search fees were timely paid by the applicant, this International Search Report covers only those claims for which fees were paid, specifically claims Nos.:

4. ☐ No required additional search fees were timely paid by the applicant. Consequently, this International Search Report is restricted to the invention first mentioned in the claims; it is covered by claims Nos.:

Remark on Protest

- ☒ The additional search fees were accompanied by the applicant's protest.
- ☐ No protest accompanied the payment of additional search fees.

INTERNATIONAL SEARCH REPORT

International Application No. PCT/US 98/02357

FURTHER INFORMATION CONTINUED FROM PCT/ISA/ 210

This International Searching Authority found multiple (groups of) inventions in this international application, as follows:

1. Claims: 1-21

A queing system and method, where entries inputted at a network switch are placed into a plurality of queues according to the priority of each entry, where the number of priority levels matches the number of queue read sides, where the priority of an entry is included in each entry in a number of priority bits, and where the queue write side has at least twice the queueing capacity of the individual queue read sides.

2. Claims: 21-24

A queing system and method, where entries inputted at a network switch are placed into a plurality of queues according to the priority of each entry, and where an entry is sent from a queue write side to an overflow area at queue read side in case the corresponding queue read side is not capable of receiving said entry.

INTERNATIONAL SEARCH REPORT

Information on patent family members

Internati Application No

PCT/US 98/02357

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
US 4630254 A	16-12-1986	DE 3586280 A EP 0179550 A JP 1963045 C JP 6095691 B JP 61161046 A	06-08-1992 30-04-1986 25-08-1995 24-11-1994 21-07-1986
EP 0544188 A	02-06-1993	JP 5153152 A CA 2082318 A DE 69224461 D DE 69224461 T	18-06-1993 26-05-1993 26-03-1998 20-08-1998
US 5541922 A	30-07-1996	FI 91695 B AU 676199 B AU 4821893 A EP 0788692 A WO 9408411 A JP 8503821 T NZ 256039 A	15-04-1994 06-03-1997 26-04-1994 13-08-1997 14-04-1994 23-04-1996 26-09-1995
US 5774453 A	30-06-1998	JP 8288953 A AU 5072696 A	01-11-1996 31-10-1996
EP 0491489 A	24-06-1992	US 5166930 A AU 632006 B AU 8830391 A CA 2054355 A,C JP 2003306 C JP 4315337 A JP 7020124 B	24-11-1992 10-12-1992 18-06-1992 18-06-1992 20-12-1995 06-11-1992 06-03-1995
EP 0526104 A	03-02-1993	US 5233606 A AU 642091 B AU 2044392 A CA 2068976 A CA 2068976 C JP 2622055 B JP 7226770 A	03-08-1993 07-10-1993 11-02-1993 03-02-1993 27-09-1994 18-06-1997 22-08-1995
US 5079762 A	07-01-1992	JP 2531275 B JP 3117137 A	04-09-1996 17-05-1991

INTERNATIONAL SEARCH REPORT

Information on patent family members

Internati Application No
PCT/US 98/02357

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
US 5583858 A	10-12-1996	JP 7107093 A EP 0652685 A	21-04-1995 10-05-1995
US 5313579 A	17-05-1994	NONE	
US 5555265 A	10-09-1996	JP 7240752 A	12-09-1995